

Implementing Quality of Service

This chapter's objectives are to define the options and technical implementations for the various types of quality of service (QoS) required by enterprises for typical virtual private network (VPN) deployments. Service providers and enterprises typically build parallel networks to support the transport of data, voice, video, and mission-critical and non-mission-critical applications. With the move toward convergence, as well as the use of packet-based IP networks, the shift from circuit-switched division and parallel builds of network resources toward a single IP network is increasing. This chapter covers the requirements for supporting the converged world of Multiprotocol Label Switching (MPLS) VPNs and how this maps to QoS policy applicable in the enterprise. The aim is to provide a deployable set of policies that the enterprise can use as guidelines. You'll also see how to address these policies to the service provider. Specifically, the QoS needs of Acme, Inc. are addressed in the case study.

Introduction to QoS

Although the amount of bandwidth is increasing as higher-speed networks become more economically viable, QoS is not unnecessary. All networks have congestion points where data packets can be dropped, such as WAN links where a larger link feeds data into a smaller link, or a place where several links are aggregated into fewer trunks. QoS is not a substitute for bandwidth, nor does it create bandwidth. QoS lets network administrators control when and how data is dropped when congestion does occur. As such, QoS is an important tool that should be enabled, along with adding bandwidth, as part of a coordinated capacity-planning process.

Another important aspect to consider alongside QoS is traffic engineering (TE). TE is the process of selecting the paths that traffic will transit through the network. TE can be used to accomplish a number of goals. For example, a customer or service provider could traffic-engineer its network to ensure that none of the links or routers in the network are overutilized or underutilized. Alternatively, a service provider or customer could use TE to control the path taken by voice packets to ensure appropriate levels of delay, jitter, and packet loss.

End-to-end QoS should be considered a prerequisite with the convergence of latency-sensitive traffic, such as voice and videoconferencing along with more traditional IP data

traffic in the network. QoS becomes a key element in delivery of service in an assured, robust, and highly efficient manner. Voice and video require network services with low latency, minimal jitter, and minimal packet loss. The biggest impact on this and other real-time applications is packet loss and delay, which seriously affects the quality of the voice call or the video image. These and other data applications also require segregation to ensure proper treatment in this converged infrastructure.

The application of QoS is a viable and necessary methodology to provide optimal performance for a variety of applications in what is ultimately an environment with finite resources. A well-designed QoS plan conditions the network to give access to the right amount of network resources needed by applications using the network, whether they are real-time or noninteractive applications.

Before QoS can be deployed, the administrator must consider developing a QoS policy. Voice traffic needs to be kept separate because it is especially sensitive to delay. Video traffic is also delay-sensitive and is often so bandwidth-intensive that care needs to be taken to make sure that it doesn't overwhelm low-bandwidth WAN links. After these applications are identified, traffic needs to be marked in a reliable way to make sure that it is given the correct classification and QoS treatment within the network.

If you look at the available options, you must ask yourself some questions that will inevitably help guide you as you formulate a QoS strategy:

- Do I need to support real-time delay-sensitive applications?
- Do I have mission-critical applications that require special handling?
- Do I know which applications and services are being planned that may affect the strategy?
- Does my selection correspond with what I am being offered by the service provider? If not, how do I make this transparent?
- What current traffic patterns or aggregate application traffic should I take into consideration?

From these questions, you have various options: Define a policy that supports the use of real-time applications and that treats everything else as best-effort traffic, or build a tiered policy that addresses the whole. After all, QoS can provide a more granular approach to segmentation of traffic and can expedite traffic of a specific type when required. You will explore these options in this chapter.

After the QoS policies are determined, you need to define the “trusted edge,” which is the place where traffic is marked in a trustworthy way. It would be useless to take special care in transporting different classes of network traffic if traffic markings could be accidentally or maliciously changed. You should also consider how to handle admission control—metered access to finite network resources. For example, a user who fires up an application that consumes an entire pipe and consequently affects others' ability to share the resource needs a form of policing.

Building a QoS Policy: Framework Considerations

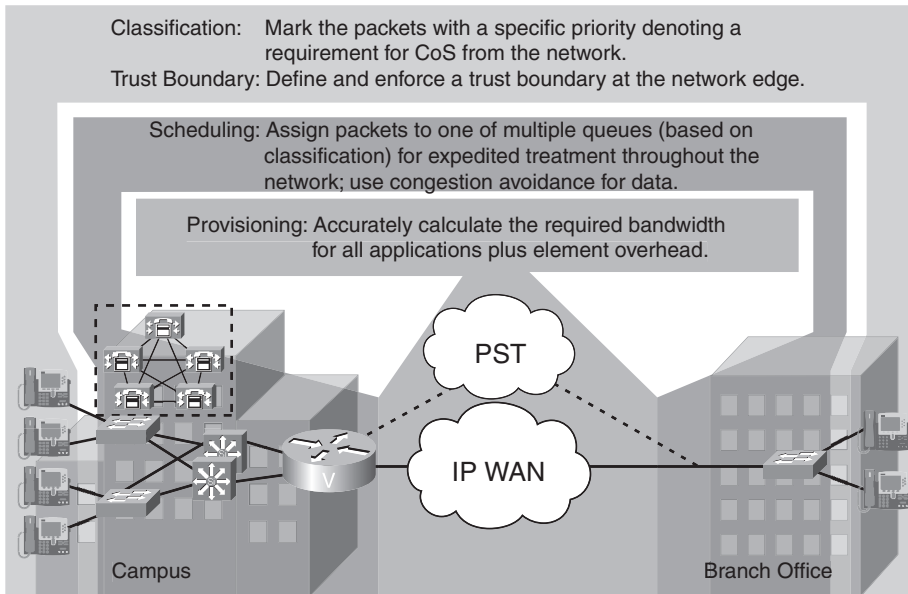
Traffic on a network is made up of flows, which are placed on the wire by various functions or endpoints. Traffic may consist of applications such as Service Advertising Protocol (SAP), CAD/CAM, e-mail, voice, video, server replication, collaboration applications, factory control applications, branch applications, and control and systems management traffic.

If you take a closer look at these applications, it is apparent that some level of control over performance measures is necessary—specifically, the bandwidth and delay/jitter and loss that each class of application can tolerate. These performance measures can vary greatly and have various effects. If you apply a service level against these performance measures, it can be broadly positioned into four levels that drive the strategy:

- **Provisioning**—The first step is ensuring that the correct transport is selected. Appropriate allocation of bandwidth ensures the proper start point for network design. Understanding application characteristics is key—what they will use in terms of network bandwidth and their delay, jitter, latency, and loss needs.
- **Best-effort service**—The majority of application data flows fit this service level. Best-effort service provides basic connectivity with no guarantee for packet delivery and handling.
- **Differentiated service**—Traffic at this service level can be grouped into classes based on their individual requirements. Each class is then treated according to its configured QoS mechanism.
- **Guaranteed service**—Guaranteed service requires absolute allocation of specific resources to ensure that the traffic profiled to receive this service has its specific requirements met.

Applying these service levels against the application classes for their required level of service means that you need to understand where in the network they should be applied. The best approach is to define a “trust boundary” at the edge of the network where the endpoints are connected, as well as look at the tiers within the network where congestion may be encountered. After you know these, you can decide on the policy of application.

For example, in the core of the network, where bandwidth may be plentiful, the policy becomes a queue scheduling tool. However, at the edge of the network, especially where geographically remote sites may have scarce bandwidth, the policy becomes one of controlling admission to bandwidth. Basically, this is equivalent to shoving a watermelon down a garden hose—intact! Figure 5-1 outlines the high-level principles of a QoS application in network design.

Figure 5-1 *Principles of QoS Application*

This design approach introduces key notions to the correct road to QoS adoption. These notions provide the correct approach to provisioning before looking at classification of packets toward a requirement for a class of service (CoS) over the network. Determine where the trust boundary will be most effective before starting such a classification, and then indicate the area of the network where scheduling of packets to queues is carried out. Finally, determine the requirement of provisioning that is needed to ensure that sufficient bandwidth exists to carry traffic and its associated overheads.

After the network's QoS requirements have been defined, an appropriate service model must be selected. A service model is a general approach or a design philosophy for handling the competing streams of traffic within a network. You can choose from four service models:

- Provisioning
- Best-effort
- Differentiated Services (DiffServ)
- Guaranteed Services or Integrated Services (IntServ)

Provisioning is quite straightforward. It is about ensuring that there is sufficient base capacity to transport current applications, with forward consideration and thinking about future growth needs. This needs to be applied across the LANs, WANs, and MANs that will

support the enterprise. Without proper consideration to provisioning appropriate bandwidth, QoS is a wasted exercise.

The best-effort model is relatively simple to understand because there is no prioritization and all traffic gets treated equally regardless of its type. The two predominant architectures for QoS are DiffServ, defined in RFC 2474 and RFC 2475, and IntServ, documented in RFC 1633, RFC 2212, and RFC 2215. In addition, a number of RFCs and Internet Drafts expand on the base RFCs—particularly RFC 2210, which explores the use of RSVP with IntServ. Unfortunately, the IntServ/RSVP architecture does not scale in large enterprises due to the need for end-to-end path setup and reservation. The service model selected must be able to meet the network's QoS requirements as well as integrate any networked applications. This chapter explores the service models available so that you can leverage the best of all three.

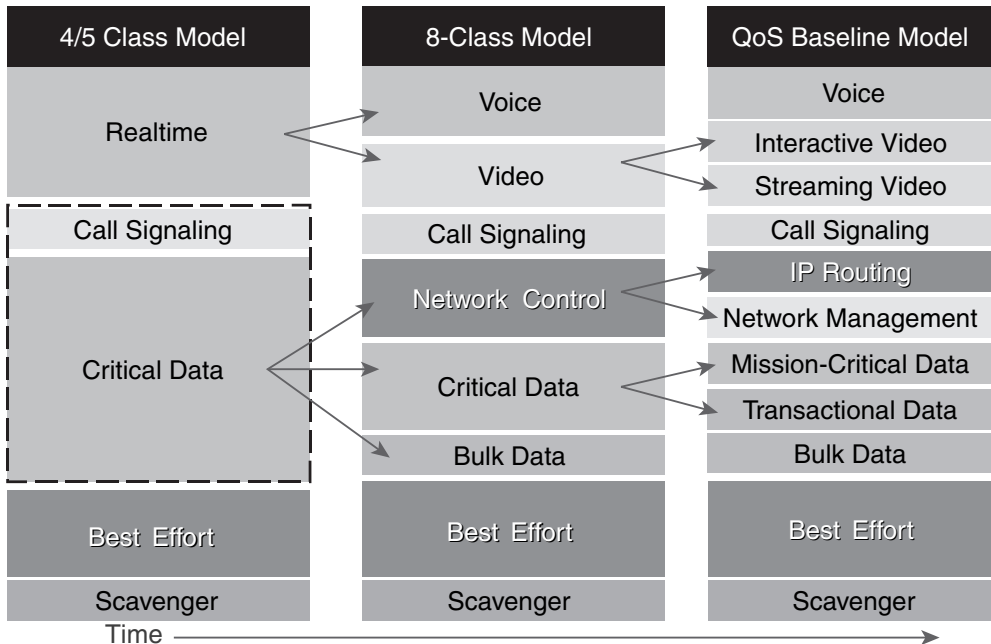
Implementing QoS is a means to use bandwidth efficiently, but it is not a blanket substitute for bandwidth itself. When an enterprise is faced with ever-increasing congestion, a certain point is reached where QoS alone does not solve bandwidth requirements. At such a point, nothing short of another form of QoS or correctly sized bandwidth will suffice.

QoS Tool Chest: Understanding the Mechanisms

The preceding section discussed the need to apply a QoS policy in the network, where it should be applied, and the models under which you can operate. To help you completely understand the mechanics of QoS, the next section explores the mechanisms available to perform the tasks at hand.

Classes of Service

To provide a mechanism for prioritizing the different types of IP traffic that exist on the network, it is important to adopt a CoS model that is flexible and simple to maintain and that meets the behavioral needs of different applications. Applications can then be categorized into the appropriate classes according to their delivery requirements. Based on this strategy, the following QoS classes of service are defined to address the different forwarding requirements of all traffic while maintaining a small number of classes. Figure 5-2 shows an approach an enterprise could follow, leading toward the mature 11-class IP precedence model. A four- or five-class IP precedence model should be the starting baseline an enterprise should consider. This model allows a migration path as more granular classes are added over time.

Figure 5-2 *How Many Classes of Service Do You Need?*

Understanding the deployment needs is a multistep process:

- Step 1** Strategically define the business objectives to be achieved via QoS.
- Step 2** Analyze the service-level requirements of the various traffic classes to be provisioned for.
- Step 3** Design and test QoS policies before production network rollout.
- Step 4** Roll out the tested QoS designs to the production network.
- Step 5** Monitor service levels to ensure that the QoS objectives are being met.

These steps may need to be repeated as business conditions change and evolve. These steps are derived from the QoS baseline model developed by Tim Sziget.

The classifications are split into two different areas: Layer 3 classification and Layer 2 CoS. The Layer 3 classifications cover the following:

- IP Precedence (or type of service [ToS]) markings
- Differentiated Services Code Point (DSCP), which provides for markings based on value ranges, where each DSCP specifies a particular per-hop behavior that is applied to a packet

- Per-hop behavior forwarding treatment applied at a differentiated services-compliant node to a behavior aggregate

IP ToS

In general, when referring to the ToS values, there are two methods for specifying QoS information within an IP packet. The first is to use the three most-significant bits (MSBs) of the ToS field in the IP header. These are called the IP Precedence (IPP) values. They allow for up to eight user-definable classes of service. The second method, referring to the 6 MSBs of the ToS field, is an extension of the IPP model, which allows for up to 64 DSCP values.

Based on these classifications, real-time voice bearer traffic is marked as Class 5 with guaranteed expedited delivery using an expedited queuing mechanism for voice traffic to ensure that voice quality is not adversely affected under heavy link utilization. This mechanism alone cannot guarantee protection for voice, so it needs to be used in combination with good capacity planning and call admission control (CAC). You will explore the capabilities available in the upcoming sections.

Traffic marked as Class 2, 3, 4, and 6 is provided guaranteed minimum bandwidth and is serviced via class-based weighted fair queuing (CBWFQ). The minimum bandwidth used should be calculated to account for peak usage for all traffic within each class. Should these classes require bandwidth usage that exceeds the configured minimum amount, this would be allowed, provided that other classes are not fully using their minimum bandwidth allocation.

All traffic marked as Class 0 is guaranteed the remainder of the bandwidth. Class 1 (batch/scavenger) traffic is drop-insensitive, or batch transfers are given a lower-priority treatment than all other classes. Typically, Class 1 should be assigned the smallest possible amount of bandwidth. Therefore, in the event of link congestion, Class 1's bandwidth usage is immediately contained to protect other higher-priority data.

Although the standard direction is to move toward the full adoption of the DSCP model, older implementations used to define the classes of service and perform traffic matching and queuing based on the IPP values. Because the IPP value is based on the first 3 MSBs of the DSCP field, it is possible for each IPP value to cover the full range of DSCP drop precedence values (bits 3 to 6) for each class selector. It should be noted that such mechanisms are now better placed to be moved to DSCP support to allow for the additional benefit of expedited forwarding/assured forwarding (EF/AF) class granularity and scaling of classes supported over time.

Ensure that the correct traffic mapping is carried out. Failing to do so may lead to classification of voice traffic to some value other than the DSCP value of 46 (EF). This may come about as a result of classification errors within the network or at the LAN edge due to incorrect CoS-to-DSCP or DSCP-to-CoS mappings, which can lead to service impact.

Hardware Queuing

QoS-enabled Ethernet switches provide a Layer 2 (L2) queuing mechanism, which allows for ingress and egress queuing. Ingress frames arriving at the L2 switch require buffering before scheduling on the egress port. Therefore, depending on the number of buffers available to each port, it is possible for ingress frames to be dropped instantaneously. If strict priority queues are not used, real-time voice traffic is not guaranteed for expedited delivery. Using the priority queues, if present, for both ingress and egress traffic provides a low-latency path through the L2 device for delay-sensitive traffic. All current Cisco platforms (2950, 2970, 3550, 3750, 4500, and 6500) support the use of internal DSCP to determine QoS treatment. These are derived from either the packets' DSCP trust classification, the trusted CoS markings, or an explicit configuration policy.

Although no standard number of queues is provided, the port capabilities can be determined via Cisco IOS or CatOS. The information is presented separately for both transmit and receive interfaces and is represented in 1P x Q y T format. 1P refers to the strict priority queue available, x Q refers to the number of input or output queues available, and y T is the number of drop or Weighted Random Early Detection (WRED) thresholds that can be configured. It is recommended that all future hardware support a minimum of 1P1Q queuing for both ingress and egress.

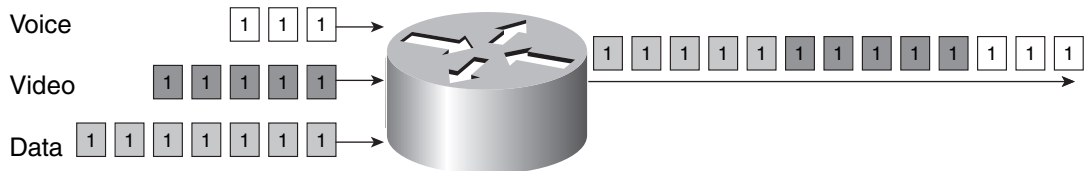
Software Queuing

Prioritization and treatment of traffic is based on defined CoSs. Where potential network congestion may occur ensures that each class receives the appropriate forwarding priority, as well as minimum reserved bandwidth. If traffic for each class exceeds the allocated bandwidth requirements, depending on the WAN technologies, one of the following actions needs to be taken:

- Drop the excess traffic
- Forward excess traffic without changes to the original QoS information
- Forward the excess traffic with the ToS bits reset to a lower-priority CoS

QoS Mechanisms Defined

The QoS architecture introduces multiple components that form the basis of the building blocks of an end-to-end solution. First, you must understand the various capabilities that are available, as shown in Figure 5-3.

Figure 5-3 *Scheduling Tools: Queuing Algorithms*

These capabilities can be broken into several categories:

- **Classification and marking**—Packet classification features allow traffic to be partitioned into multiple priority levels or CoSs.

Packets can be classified based on the incoming interface, source or destination addresses, IP protocol type and port, application type (network-based application recognition [NBAR]), IPP or DSCP value, 802.1p priority, MPLS EXP field, and other criteria. Marking is the QoS feature component that “colors” a packet (frame) so that it can be identified and distinguished from other packets (frames) in QoS treatment. Policies can then be associated with these classes to perform traffic shaping, rate-limiting/policing, priority transmission, and other operations to achieve the desired end-to-end QoS for the particular application or class. Figure 5-2 showed an overview of classification for CoS, ToS, and DSCP.

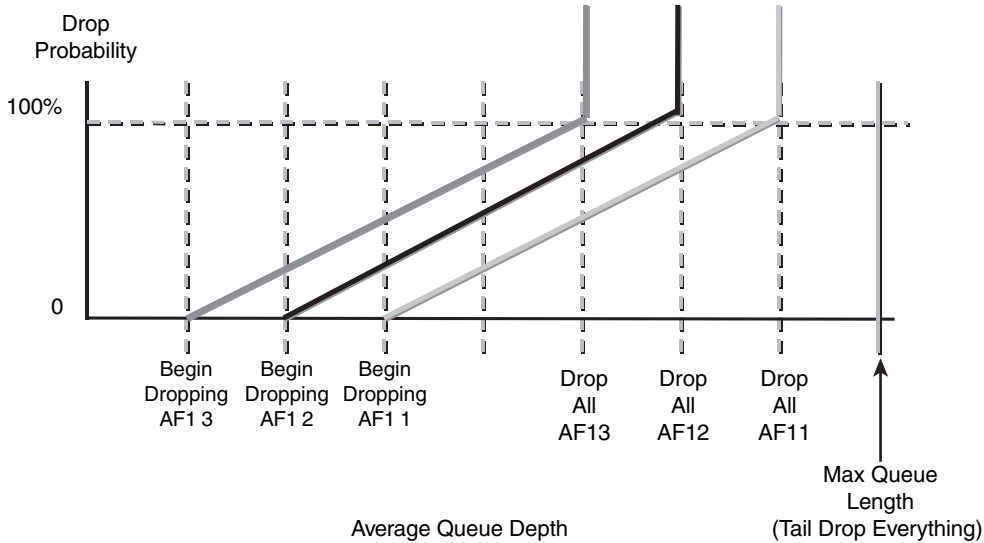
- **Congestion management**—Congestion-management features control congestion after it occurs.

Queuing algorithms are used to sort the traffic and then determine some method of prioritizing it onto an output link. Congestion-management techniques include Weighted Fair Queuing (WFQ), CBWFQ, and low-latency queuing (LLQ):

- WFQ is a flow-based queuing algorithm that does two things simultaneously: It schedules interactive traffic to the front of the queue to reduce response time, and it fairly shares the remaining bandwidth between high-bandwidth flows.
- CBWFQ guarantees bandwidth to data applications.
- LLQ is used for the highest-priority traffic, which is especially suited for voice over IP (VoIP).
- **Congestion avoidance**—Congestion-avoidance techniques monitor network traffic loads in an effort to anticipate and avoid congestion at common network and internetwork bottlenecks before it becomes a problem.

As shown in Figure 5-4, the WRED algorithm avoids congestion and controls latency at a coarse level by establishing control over buffer depths on both low- and high-speed data links. WRED is primarily designed to work with TCP applications. When WRED is used and the TCP source detects the dropped packet, the source slows its transmission. WRED can selectively discard lower-priority traffic when the interface begins to get congested.

Figure 5-4 Congestion Avoidance: DSCP-Based WRED



- **Traffic conditioning**—Traffic entering a network can be conditioned (operated on for QoS purposes) by using a policer or shaper.

Traffic shaping involves smoothing traffic to a specified rate through the use of buffers. A policer, on the other hand, does not smooth or buffer traffic. It simply re-marks (IPP/DSCP), transmits, or drops the packets, depending on the configured policy. Legacy tools such as committed access rate (CAR) let network operators define bandwidth limits and specify actions to perform when traffic conforms to, exceeds, or completely violates the rate limits. Generic traffic shaping (GTS) provides a mechanism to control traffic by buffering it and transmitting at a specified rate. Frame Relay traffic shaping (FRTS) provides mechanisms for shaping traffic based on Frame Relay service parameters such as the committed information rate (CIR) and the backward explicit congestion notification (BECN) provided by the Frame Relay switch.

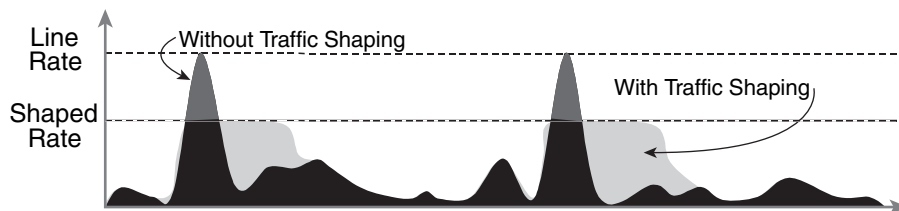
Policers and shapers are the oldest forms of QoS mechanisms. These tools have the same objectives—to identify and respond to traffic violations. Policers and shapers usually identify traffic violations in an identical manner; however, their main difference is the manner in which they respond to violations:

- A policer typically drops traffic.
- A shaper typically delays excess traffic using a buffer to hold packets and shape the flow when the source's data rate is higher than expected.

The principal drawback of strict traffic policing is that TCP retransmits dropped packets and throttles flows up and down until all the data is sent (or the connection times out). Such TCP ramping behavior results in inefficient use of bandwidth, both overutilizing and underutilizing the WAN links.

Since shaping (usually) delays packets rather than dropping them, it smoothes flows and allows for more efficient use of expensive WAN bandwidth. Therefore, shaping is more suitable in the WAN than policing. Figure 5-5 demonstrates the need for policers and shapers.

Figure 5-5 Provisioning Tools: Policers and Shapers



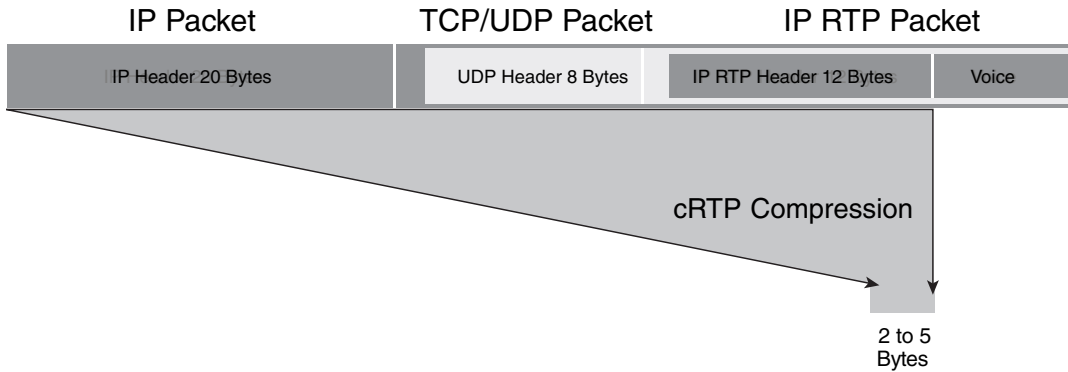
This is especially the case with nonbroadcast multiaccess (NBMA) WAN media, where physical access speed can vary between two endpoints, such as Frame Relay and ATM.

- **Link efficiency mechanisms**—Two link efficiency mechanisms work in conjunction with other QoS features to maximize bandwidth utilization.

Newer multimedia application traffic, such as packetized audio and video, is in Real-Time Transport Protocol (RTP) packets and Cisco IOS Software. This saves on link bandwidth by compressing the RTP header (Compressed Real-Time Protocol [cRTP]), as shown in Figure 5-6.

VoIP packets are relatively small, often with a G.729 voice packet and approximately 20-byte payloads. However, the IP plus User Datagram Protocol (UDP) plus RTP headers equal 40 bytes (uncompressed), which could therefore account for nearly two-thirds of the entire packet. A solution is to use the Van Jacobsen algorithm to compress the headers for VoIP, reducing the header size from 40 bytes to less than 5 bytes.

Figure 5-6 *IP RTP Header Compression*

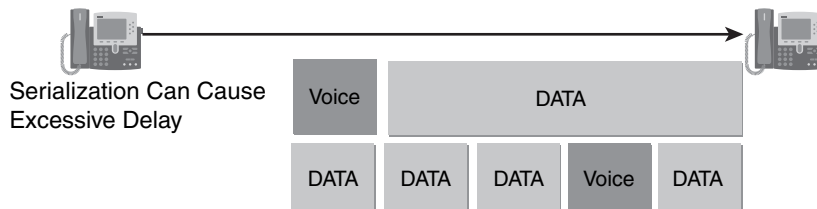


A data frame can be sent to the physical wire at only the interface’s serialization rate. This serialization rate is the frame’s size divided by the interface’s clocking speed. For example, a 1500-byte frame takes 214 ms to serialize on a 56-kbps circuit.

If a delay-sensitive voice packet is behind a large data packet in the egress interface queue, the end-to-end delay budget of 150 ms could be exceeded. Refer to ITU G.114, which defines this measure as “most users being satisfied.” Additionally, even a relatively small frame can adversely affect overall voice quality by simply increasing the jitter to a value greater than the size of the adaptive jitter buffer at the receiver.

Link Fragmentation and Interleaving (LFI) tools fragment large data frames into regular-sized pieces and interleave voice frames into the flow so that the end-to-end delay can be accurately predicted. This places bounds on jitter by preventing voice traffic from being delayed behind large data frames. To decrease latency and jitter for interactive traffic, LFI breaks up large datagrams and interleaves delay-sensitive interactive traffic with the resulting smaller packets, as shown in Figure 5-7.

Figure 5-7 *LFI*



A maximum of 10-ms serialization delay is the recommended target to use for setting fragmentation size. This allows for headroom on a per-hop basis, because it allows adequate time for end-to-end latency required by voice. Two tools are available for LFI: multilink PPP (MLP) LFI (for point-to-point links) and FRF.12 (for Frame Relay links).

While reviewing these capabilities, it is important to keep in mind that the LLQ is in effect a first-in, first-out (FIFO) queue. The amount of bandwidth reserved for the LLQ is variable, yet if the LLQ is overprovisioned, the overall effect is a dampening of QoS functionality. This is because the scheduling algorithm that decides how packets exit the device is predominantly FIFO. Overprovisioning the LLQ defeats the purpose of enabling QoS. For this reason, it is recommended that you not provision more than 33 percent of the link's capacity as LLQ.

The 33 percent limit for all LLQs is a design guideline recommendation only. There may be cases where specific business needs cannot be met while holding to this recommendation. In such cases, the enterprise must provision queuing according to its specific requirements and constraints.

To avoid bandwidth starvation of background applications such as network management services and best-effort traffic types, it is recommended that you not provision total bandwidth guarantees to exceed 75 percent of the link's capacity. This is a subjective area because of the size of the link employed for the connection. It should be applied as a general rule. When a larger link is employed, the 75 percent rule is less relevant—links such as E3/DS3 and above, for example. Figure 5-8 provides an overview of the recommendations for WAN egress design scheduling.

Figure 5-8 WAN Scheduling Design Principles

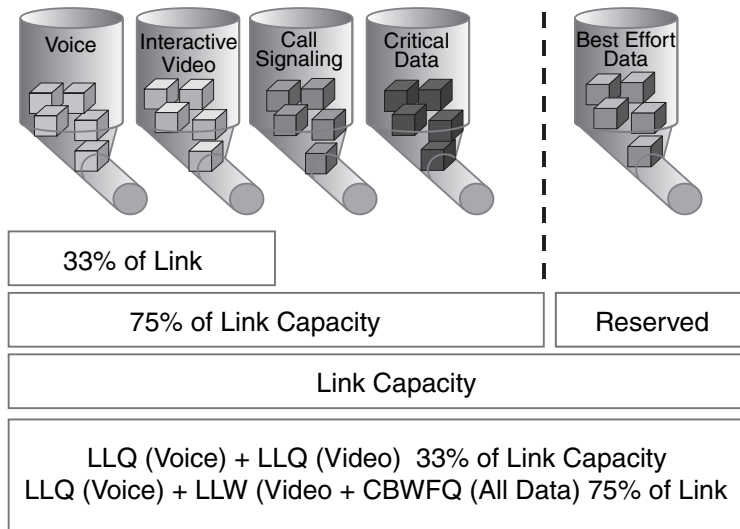
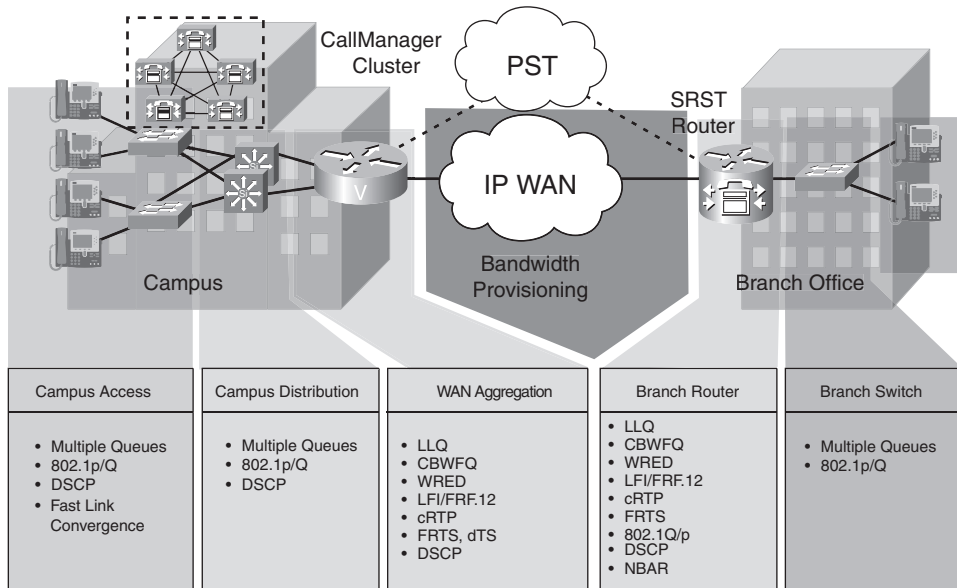


Figure 5-9 blends this all together by demonstrating the relevant tools and how they apply in the context of the network. You will explore how and where these should be applied in the case study later in this chapter.

Figure 5-9 *QoS Tools Mapping*



From the tools mapping, you can see that you have various options to use in terms of campus access, distribution, and WAN aggregation, with capabilities that extend through the service provider MPLS VPN and into the branch network.

Pulling It Together: Build the Trust

As discussed earlier, there are many places in the network in which the application of QoS, either marking or classification, occurs. In this section, you will pull this together in some configurations, starting with the trust boundary principle of marking nearest to the end-points on the network.

To apply this, you need to understand the edge marking mechanisms that can be applied. In this case, you will use the notion of trust boundaries. The concept of trust is an important and integral one to implementing QoS. As soon as the end devices have a set CoS or ToS, the switch can either trust them or not. If the device at the edge (in our case, a switch) trusts the settings, it does not need to do any reclassification. If it does not trust the settings, it must perform reclassification for the appropriate QoS.

The notion of trusting or not trusting forms the basis of the trust boundary. Ideally, classification should be done as close to the source as possible. If the end device can perform this function, the trust boundary for the network is at the access layer. This depends on the capabilities of the switch in the access layer. If the switch can reclassify the packets, the trust boundary remains in the access layer. If the switch cannot perform this function, the task falls to other devices in the network going toward the backbone.

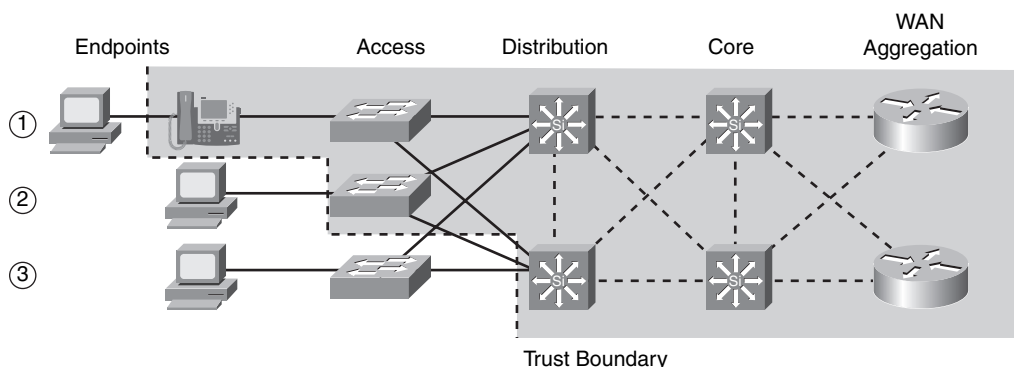
In this case, the rule of thumb is to perform reclassification at the distribution layer. This means that the trust boundary has shifted to the distribution layer. It is more than likely that there is a high-end switch in the distribution layer with features to support this function. If possible, try to avoid performing this function in the core of the network.

Frames and packets can be marked as important by using Layer 2 CoS settings in the User Priority bits of the 802.1p portion of the 802.1Q header or the IPP/DSCP bits in the ToS byte of the IPv4 header.

Figure 5-10 gives an overview of the trust boundary states that can be applicable when establishing this capability:

- A device is *trusted* if it correctly classifies packets.
- For scalability, classification should be done as close to the edge as possible.
- The outermost trusted devices represent the trust boundary.
- (1) and (2) are optimal; (3) is acceptable (if the access switch cannot perform classification).

Figure 5-10 Establishing the Trust Boundary



For example, suppose you have a LAN edge that is configured to have the voice traffic sit in an auxiliary virtual LAN (VLAN) and data traffic that transports standard desktop PC connectivity in a standard data VLAN. In this case, you can establish a policy of trust on the auxiliary VLAN where the voice endpoints are connected and there's no trust for the data VLAN. This forms a fundamental design principle of the Differentiated Services model, which is to classify and mark packets as close to the source as possible. To keep

users from marking their own traffic, a trust boundary needs to be enforced, which should likewise be as close to the source as possible.

Building the Policy Framework

To achieve a consistent end-to-end network QoS implementation, it is expected that all IP traffic follow a set of defined processes from the source to the destination device. As a strategy, the trust boundary should be established at the LAN edge closest to the connected devices, such as desktop/voice/server L2 switches or lab gateways. All IP traffic that arrives at the LAN edge switch should be classified as either trusted or untrusted.

Classification and Marking of Traffic

As the name implies, trusted devices, such as IP phones, call managers, and unity voice servers, already originate traffic with the desired ToS marked to the appropriate values. Hence, it is simply a matter of trusting and preserving the received ToS values when the packets are subsequently forwarded to the next switch or router in the network. On the other hand, not all traffic originating from devices that are user-/admin-configurable, such as desktop PCs and file/print/application servers, should be trusted with respect to their ToS settings (real-time desktop applications are covered in later sections). Therefore, the traffic from these devices needs to be re-marked with the appropriate ToS values that accurately reflect the desired level of priority for such traffic up to a predefined bandwidth limit.

As soon as the traffic is classified as trusted or is re-marked with the correct ToS settings at the LAN edge into one of the defined classes of service, a trusted edge boundary is established. This enables the traffic to be fully trusted within the network such that it can be prioritized and acted on accordingly at any of the potential congestion points. Typically, these congestion points are at the WAN edge; however, they can also be found at the LAN aggregations.

Trusted Edge

For traffic to be fully trusted within the network core, it is critical that all traffic classification and re-marking at the edge be performed before any frames are forwarded to the next switch or router. This means that the ingress traffic must be inspected to determine whether it is to be trusted or untrusted. If it is considered trusted, the L2 CoS, Layer 3 (L3) IPP, or DSCP values can be derived from the incoming frame/packet and subsequently forwarded to the next device unmodified. Untrusted traffic, on the other hand, should be rewritten to a default DSCP value.

Device Trust

To simplify the edge classification and re-marking operation, the concept of a trusted device needs to be defined. These are devices that are known to correctly provide QoS markings for traffic they originate. Furthermore, these devices also have limited or minimal user QoS configuration capability, such as IP phones, Call Manager/Unity/IPCC servers, and voice/videoconferencing gateways. Whenever these devices are connected to the L2 edge switch, it is easier to trust the L2 CoS or L3 DSCP information on the ingress port rather than to manually identify the type of traffic that should be trusted by means of access control lists (ACLs). For some traffic types, such as RTP streams, it is more difficult to match on specific Layer 4 (L4) ports because the application operates on dynamic port ranges. In such cases, where possible, it is preferable to allow the device or application to correctly classify the traffic and be trusted by the switch.

Application Trust

Although it is possible to establish a trust boundary using ingress CoS or ToS values from devices that are considered trusted, it is important to note that not all devices support the proper QoS marking. Hence, forwarding the traffic without first modifying the ToS value to the appropriate IPP or DSCP can potentially result in erroneous QoS treatment of the traffic for a particular CoS at the WAN edge. By passing all untrusted traffic through an ACL at the LAN edge, it is possible to correctly identify applications that cannot provide correct QoS marking based on the L3/L4 protocol information. Subsequently, these applications can be reclassified and marked to the appropriate classes of service. All other traffic that the ACL does not correctly identify should have its CoS and/or ToS values rewritten to the default/best-effort CoS.

An example of classification by ACL is to re-mark traffic originating Cisco Softphone RTP stream from workstations to Class 5 and associated Skinny (SCCP) packets to Class 3 and non-drop-sensitive batch transfer traffic to Class 1. All other traffic is rewritten to Class 0 regardless of how the original ToS values are set.

By enabling trust for specific real-time desktop applications, such as Cisco Softphone and videoconferencing, it is envisioned that a strategy for ingress traffic policing or rate limiting of traffic belonging to Classes 3, 4, and 5 also be applied at the LAN switch. This would ensure that each attached desktop machine does not exceed a predefined maximum bandwidth value for these priority classes of service. This does not eliminate the need for a comprehensive call admission control (CAC) implementation for voice and video. However, this is yet another level of protection against potential misuse of the QoS classes of service whether it is executed intentionally or unintentionally.

CoS and DSCP

At the L2 edge, Cisco IP phones can VLAN trunk to the Catalyst switches using 802.1Q tagging. Within the 802.1Q tagged frame is a 3-bit CoS field that is commonly referred to as the 802.1p bits. Coincidentally, this is equivalent to the 3 IPP bits within the L3 ToS field. Hence, to maintain end-to-end QoS, it is necessary to ensure that the CoS-to-IPP and IPP-to-CoS are consistently mapped throughout the network. Similarly, CoS values can also be mapped to DSCP to provide the same end-to-end QoS functionality; however, care must be taken to ensure that each CoS value is mapped to a DSCP range of values that has a common 3 MSBs.

By leveraging the CoS setting of the frame coming from the phone, strict priority ingress queuing is possible on Catalyst platforms that support a receive queue mechanism. Because the intelligence of the ingress application-specific integrated circuit (ASIC) on the switch is limited to L2 header inspection, ingress queuing based on L3 ToS values is not possible. For trusted devices that are not capable of 802.1Q trunking to the Ethernet switches, as well as ports configured for edge switch to router uplinks, it is necessary to trust the DSCP values of all incoming IP packets.

Strategy for Classifying Voice Bearer Traffic

Voice traffic traverses the network in the form of RTP streams. The Cisco IP phone originates RTP packets with a DSCP value of 46 (EF) and a CoS value of 5. Based on the device trust model, the DSCP value of voice packets should be preserved across the network. Because MPLS VPN is the WAN transport you are considering, it is common across many service providers to expect RTP traffic presented at the ingress of the provider edge (PE) device to be marked with a DSCP value of 46.

QoS on Backup WAN Connections

Today, WAN services are diverse. Depending on the geographic location of the sites, they may include technologies such as ATM, Frame Relay, ISDN, point-to-point time-division multiplexing (TDM), and network-based VPNs. However, a site is not always provisioned with equal-sized connections. Hence, if a backup connection exists, it is expected to be of a lower bandwidth than the primary link as well as being idle during normal operating conditions. This means that when the site's traffic is required to be routed over the backup connection, potential oversubscription of the link may occur.

To understand how QoS can be applied to back up WAN circuits, it is important to understand exactly how much is allocated for each CoS on the primary connection. However, due to the diverse nature of the types of site locations and sizes of WAN circuits implemented in today's environment, the overall amount of bandwidth required for real-time traffic can vary from one site to another. It is therefore recommended that, for any given fixed-line primary WAN link, no more than 33 percent of the total available bandwidth be assigned to

traffic belonging to Class 5. This is consistent with the overall recommendations for provisioning LLQ, as discussed in earlier sections. This can also be overprovisioned.

Shaping/Policing Strategy

There are many cases in which a connection's guaranteed bandwidth is not necessarily the same as the physical port speed. Hence, anything that is transmitted in excess of the available bandwidth is subject to policing and potentially can be dropped by the service provider without regard for the traffic classes. WAN technologies such as Frame Relay, ATM, and L2 and L3 IP/VPN services are good examples of this, whereby it is possible to transmit up to the physical access port speed. However, the service provider guarantees delivery for traffic only up to the contracted bandwidth such as CIR or sustainable cell rate (SCR) in the case of ATM.

The decision of whether excess traffic should be dropped or marked down to a different class depending on the applicable WAN technology should be left to the discretion of the network administrator. Traffic shaping and policing at the WAN edge means that there is more control over the type of excess traffic that should be sent to the provider's network. This avoids the chance that service providers will indiscriminately discard excess traffic that belongs to all classes of service.

A better approach is to treat traffic with a different drop preference, as defined in RFC 2597, *Assured Forwarding Drop Preference*. The reason for this is that different queues are drained at different rates. Therefore, if you mark to a different class, you introduce out-of-sequence packet delivery, which has detrimental effects on the in-profile traffic. DSCP-based WRED is then employed to discard out-of-profile traffic aggressively ahead of in-profile traffic.

Real-time voice traffic is sensitive to delay and jitter. Therefore, it is recommended that whenever possible, excess burst traffic for this class should be policed. By default, real-time voice traffic that exceeds the bandwidth allocated to the strict priority queue (low-latency queuing [LLQ]) is allowed limited bursting.

When using network-based VPN services, such as MPLS VPN, and depending on the service provider offerings, each class can be allocated guaranteed bandwidth across the provider network. By shaping each class at the customer edge (CE), excess traffic could still be forwarded to the provider network but may be marked down to a different CoS or set to a higher drop preference value within the same class selector. This would ensure that if a service provider experiences network congestion, traffic is dropped based on the network administrator's preference rather than random discards by the provider.

NOTE Providers apply an always-on policer to protect their core capacity planning. This means that even if bandwidth on the link is available, if the EF traffic exceeds allocation, it is dropped within the provider network.

In the case of L2/L3 VPN, it is also possible to have many sites connected to the same service provider network with varying connection speeds. Often, the hub site is serviced by a much larger connection while remote offices are connected at significantly reduced speed. This may cause an oversubscription of the remote WAN connection due to the peer-to-peer nature of L2/L3 VPNs. Egress shaping should be considered on the CE, particularly at the hub location, to prevent this situation from occurring. However, some level of oversubscription of the remote link may still occur, due to the fact that remote-to-remote office traffic patterns are unpredictable and cannot be accounted for. For L2 VPN services that share a common broadcast domain, it is not recommended that these types of technology be adopted due to the difficulty inherent in egress traffic shaping.

Queuing/Link Efficiency Strategy

QoS mechanisms such as LLQ and CBWFQ address the prioritization of L3 traffic to the router's interface driver. However, the underlying hardware that places actual bits on the physical wire is made up of a single transmit ring buffer (TX ring) that operates in a FIFO fashion. The result is the introduction of a fixed delay (commonly called serialization delay) to the overall end-to-end transmission of a packet before it is encoded onto the wire. Depending on link speed and the packet's size, this serialization delay can be significant and can have a severe impact on voice quality.

Small voice packets that are processed via LLQ are still queued behind other packets in the TX ring buffer. In the worst-case scenario, the voice packet would have to wait for a 1500-byte packet to be transmitted first. Hence, serialization delay becomes a major factor in the overall end-to-end latency of the voice packet.

The size of the TX ring buffer represents a trade-off. If the buffer is too large, it is possible that too many data fragments may be placed in the queue before an LLQ fragment. This would result in the LLQ fragment's being delayed, causing higher latency and jitter. A buffer that's too small can keep higher-speed interfaces from failing to achieve line rate. For 2-Mbps circuits, the default depth is two fragments. This means that at the worst there is the potential for up to two low-priority data fragments to be on the TX ring when an LLQ fragment is ready to transmit.

Therefore, fragment sizes must be calculated to account for a transmit delay of up to two fragments. Based on the fragment-sized calculation, the worst-case jitter target for these low-speed links should be approximately 18 ms. MLP fragments in Cisco IOS have 11-byte L2 headers (with shared flags and long sequence numbers). Based on 11-byte overhead and

fragment-size multiples of 32, the following calculation is used to derive the fragment size for MLP:

$$\text{FragmentSize} = \frac{\text{LinkRate}(bps) * .009}{8} - 11(\text{Bytes})$$

Then, you round down to multiples of 32, as shown in Table 5-1.

Table 5-1 *Fragment Sizing*

Line Rate (in kbps)	Fragment Size (in Bytes)
128	128
256	256
384	416
512	576
768	864
1024	1152
> 1024	No fragmentation

The overhead of FRF.12 is similar and does not affect the fragment sizes used. As a result, LFI (via MLP or FRF.12) is required for link rates of 1024 kbps or less.

For the low-speed links defined in Table 5-2, LFI techniques such as MLP and Frame Relay FRF.12 can be used to reduce the effect of serialization delay on voice traffic in such an environment. LFI functions by fragmenting larger packets into smaller-sized fragments such that, for a given low-speed link, all packets can be transmitted onto the wire with a serialization delay that is acceptable to voice traffic. For low-speed ATM links, it is inadvisable to use MLP LFI over ATM virtual circuits because of the high overhead associated with encapsulation for small packets. Implementing low-speed ATM links is becoming an uncommon practice. However, if such scenarios exist and if voice is a requirement, it is recommended that the link speed be provisioned above 1024 kbps to avoid serialization delay and the need for fragmentation.

Table 5-2 LFI Fragment Sizing: Serialization Delay for Link Speeds of 1024 kbps and Below

Link Speed (in kbps)	Packet Size (in Bytes)					
	64	128	256	512	1024	1500
64	8 ms	16 ms	32 ms	64 ms	128 ms	187 ms
128	4 ms	8 ms	16 ms	32 ms	64 ms	93 ms
256	2 ms	4 ms	8 ms	16 ms	32 ms	46 ms
512	1 ms	2 ms	4 ms	8 ms	16 ms	23 ms
768	0.64 ms	1.3 ms	2.6 ms	5.1 ms	10.3 ms	15 ms
1024	0.4 ms	0.98 ms	2 ms	3.9 ms	7.8 ms	11.7 ms

For WAN services that do not support LFI, such as digital subscriber line (DSL) and cable technologies, it is recommended that manual override of TCP segment size values be configured on the connected router interface. This ensures that large TCP packets are fragmented to reduce the effect of serialization delay on low-speed broadband circuits.

IP/VPN QoS Strategy

Layer 3 VPN technology, such as MPLS VPN, introduces several challenges. One of those challenges is the QoS treatment and handling of traffic across the service provider's IP network, which would likely have a different type and number of QoS CoSs. Given that traffic would be routed by the provider across the IP network, it is imperative that the internal QoS classes of service be handled correctly to ensure that service levels are being met.

In some cases, there may not be a direct one-to-one mapping of enterprise CoSs to those offered by the service providers. In that case, it is necessary at the WAN edge to merge or remap the internal classes so that they may align. To ensure that important and high-priority classes are given the same level of service as if they were traversing the internal private WAN, care must be taken when such actions are carried out.

Enterprises implement more classes of service, because they want to separate applications. However, in the provider's IP core, they aggregate classes based on the service-level agreement (SLA) type. That is, they have priority queuing (PQ) for controlled latency and CBWFQ for guaranteed-bandwidth and best-effort. All CBWFQ applications that are separated at the edge are lumped together. However, as long as the aggregate meets the needs of the sum of the individual guarantees at the edge, it is fine for the service provider core and is of no concern.

Service providers may have different strategies for enforcing QoS classes. Although it may be a common practice for one provider to discard excess traffic marked with higher drop

precedence within a class selector, others may elect to drop traffic from lower-priority classes instead. This aspect of the provider's QoS offering must be fully understood and assessed so that the correct merging and/or remapping of an enterprise's internal classes are performed.

For example, if the service provider is offering four levels of QoS—EF, AF1, AF2, and best-effort (BE)—it is not recommended that more than one internal customer class share a common service provider class selector. That is, if traffic belonging to Class 2 is mapped to AF2, only packets that exceed the maximum bandwidth for this class should be marked down to AF22 or AF23, because these represent higher drop preference values within this particular class selector. In this case, no other traffic should be marked as AF22 or AF23, except excess traffic belonging to Class 2.

IP values 6 and 7 are also used for network control traffic (routing protocols). Most of this traffic is link-local, so an individual class of traffic can be set up for this traffic on a WAN port, with minimum bandwidth. On the CE side of the CE-to-PE links, it is recommended that a separate class be used for management traffic. On the PE side of the CE-to-PE link, this tends to vary per provider. This traffic must, at a minimum, be mapped into a high-priority data CoS in the service provider cloud.

Approaches for QoS Transparency Requirements for the Service Provider Network

Any L3 IP/VPN solution implemented in an enterprise network must support QoS transparency. QoS transparency is defined as the ability to recover your original discrete CoSs at the remote end of the IP/VPN network. It is unacceptable for multiple CoSs to be combined into one service provider class such that, at the remote end, the traffic cannot be recovered into the separate CoSs. This transparency can be achieved in one of two ways.

With the first option, the enterprise CE can convert the IP DSCP values to those expected by your service provider's PE, as long as a minimum of five discrete values across the service provider's network preserve the independence of the five CoSs. At the remote CE, traffic can be re-marked back to the appropriate levels for the enterprise network. It is unacceptable for traffic from one class to be re-marked by the network into another class such that it would end up in a different CoS when it was converted back to the enterprise's expected values.

The second option, available in MPLS VPN only, is to leave the IP DSCP markings untouched and use those values to set the MPLS Experimental (EXP) QoS bits to an appropriate level of service based on the markings defined by the enterprise.

RFC 3270 discusses more of the operational aspects with the transport of differing DiffServ implementations. It also classifies these into three effective modes. The Uniform, Pipe, and Short-Pipe modes provide the solution for service providers' flexibility in selecting how DiffServ CoSs are routed or traffic-engineered within their domain.

DiffServ tunneling modes introduce a new Per-Hop Behavior (PHB) that allows differentiated QoS in a provider's network. The tunneling mode is defined at the edge of the network, normally in the PE label switch routers (LSRs) (both ingress and egress). You may need to make changes in the P routers, and you must also consider what occurs when the topmost label is removed from a packet due to Penultimate Hop Popping (PHP). It may be necessary to copy the MPLS EXP value from the top label that is being popped to the newly exposed label; this does not always apply to all tunneling modes.

In some cases (for example, a plain non-VPN MPLS network), the PHP action on the final P router can expose a plain IP packet when a packet with only one label is received. When this IP packet is received by the egress LSR (PE), it is not possible to classify the packet based on the MPLS EXP bits because there is no label now. In these situations, you must configure the egress PE router to advertise an explicit-null label. When the PHP action is performed on the P router, a label with a value of 0 is sent. With this special label, you can mark the EXP bits as normally labeled packets, allowing the correct classification on the egress PE router.

MPLS network support of DiffServ specification defines the following tunneling modes:

- Uniform
- Pipe
- Short-Pipe

The next sections examine each tunneling mode and provide examples that show how you can configure each one. The examples include a full mapping of IPP to MPLS EXP bits. It is possible to have a number of different QoS parameters and tunneling modes for each customer.

NOTE

The configuration examples are not specific for MPLS VPN and are applicable for plain MPLS networks and Carrier-Supported Carrier (CsC) networks. It is also possible that your network can vary from another network in which many different QoS parameters and tunneling modes can be used.

Uniform Mode

DiffServ tunneling Uniform mode has only one layer of QoS, which reaches end to end. The ingress PE router (PE1) copies the DSCP from the incoming IP packet into the MPLS EXP bits of the imposed labels. As the EXP bits travel through the core, they may or may not be modified by intermediate P routers. In this example, the P router modifies the EXP bits of the top label. At the egress P router, you can copy the EXP bits to the EXP bits of the newly exposed label after the PHP. Finally, at the egress PE router, you can copy the EXP bits to the DSCP bits of the newly exposed IP packet.

Pipe Mode

DiffServ tunneling Pipe mode uses two layers of QoS:

- An underlying QoS for the data, which remains unchanged when traversing the core.
- A per-core QoS, which is separate from that of the underlying IP packets. This per-core QoS PHB remains transparent to end users.

When a packet reaches the edge of the MPLS core, the egress PE router classifies the newly exposed IP packets for outbound queuing based on the MPLS PHB from the EXP bits of the recently removed label.

Short-Pipe Mode

DiffServ tunneling Short-Pipe mode uses the same rules and techniques across the core.

The difference is that, at the egress PE router, you classify the newly exposed IP packets for outbound queuing based on the IP PHB from the DSCP value of this IP packet.

QoS CoS Requirements for the SP Network

The service provider's network must support a minimum of three classes at all interfaces with speeds of OC-12/STM-4 (622 Mbps) and less. These classes must include a real-time class using LLQ, a high-priority data class with CBWFQ "minimum bandwidth," and a best-effort class.

Four or five classes are preferred (with the minimum requirements for an LLQ class and all but one of the remaining classes supporting minimum bandwidth), such that the enterprise classes can map directly to the service provider's network.

WRED Implementations

Whereas QoS and LFI are techniques for congestion management, WRED is a technique used for congestion avoidance. WRED, when implemented, allows for the early detection of network congestion and provides the means for the selective discard of packets based on the IPP or DSCP values. When the average queue depth exceeds a user-defined minimum threshold, WRED begins discarding lower-priority packets (both TCP and UDP) based on the QoS information. The intent is to allow TCP applications to decrease their transmission rate and allow the network utilization to level out. Should the average queue depth increase above the user-defined maximum threshold, WRED reverts to "tail-drop" operation. This means that all packets entering the queue at that point are dropped until the traffic utilization is reduced to below the maximum threshold.

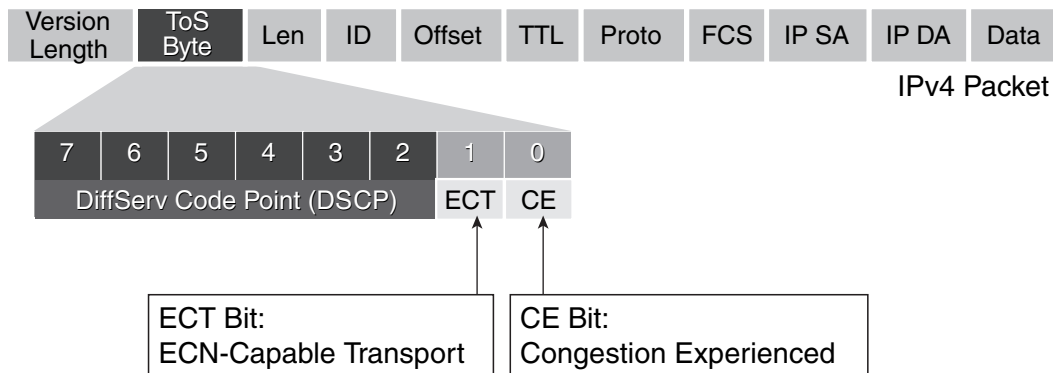
Because all traffic is classified and marked at the LAN edge, it is more useful for WRED to be implemented at the WAN edge routers. This way, when the core of the network experiences congestion, packets can be intelligently discarded. In most cases, WRED is

recommended only for WAN edge routers that directly connect to IP/VPN providers that explicitly indicate that they support this feature. Packets that exceed threshold values can have their priority marked down or selectively discarded. An important point to keep in mind is that WRED should not be applied to queues that support voice traffic, due to the potential impact that packet loss can have on voice quality.

Additionally, Explicit Congestion Notification (ECN) is an extension to WRED in that ECN marks packets instead of dropping them when the average queue length exceeds a specific threshold value. When configured with WRED's ECN feature, routers and end hosts use this marking as a signal that the network is congested and slow down sending of packets.

As stated in RFC 3168, *The Addition of Explicit Congestion Notification (ECN) to IP*, implementing ECN requires an ECN-specific field that has 2 bits: the ECN-capable Transport (ECT) bit and the CE (Congestion Experienced) bit in the IP header. The ECT bit and the CE bit can be used to make four ECN field combinations of 00 to 11. The first number is the ECT bit, and the second number is the CE bit. Figure 5-11 gives an overview of ECN application.

Figure 5-11 ECN Application



RFC3168: IP Explicit Congestion Notification

ECN is being adopted in a lot of enterprise and service provider networks, which complements WRED. The benefits can be summarized as follows:

- **Improved method for congestion avoidance**—This feature provides an improved method for congestion avoidance by allowing the network to mark packets for later transmission, rather than dropping them from the queue. Marking the packets for later transmission accommodates applications that are sensitive to delay or packet loss and provides improved throughput and application performance.

- **Enhanced queue management**—Currently, dropped packets indicate that a queue is full and that the network is experiencing congestion. When a network experiences congestion, this feature allows networks to mark a packet's IP header with a CE bit. This marking, in turn, triggers the appropriate congestion-avoidance mechanism and allows the network to better manage the data queues. With this feature, ECN-capable routers and end hosts can respond to congestion before a queue overflows and packets are dropped, providing enhanced queue management.

For more information on the benefits associated with ECN, refer to RFC 2309, *Internet Performance Recommendations*.

Identification of Traffic

After the QoS policy and toolsets that can be employed have been defined, the starting point is to understand the profile of traffic that exists in the network. To identify the traffic, several approaches can be taken to identify the current flows and start the work of classification. The easiest way to do this is to first determine the real-time applications that require special handling as they traverse the network. With these, there is a need to identify not just the bearer traffic but also the signaling traffic that may be required as part of the bearer's normal operation.

What Would Constitute This Real-Time Traffic?

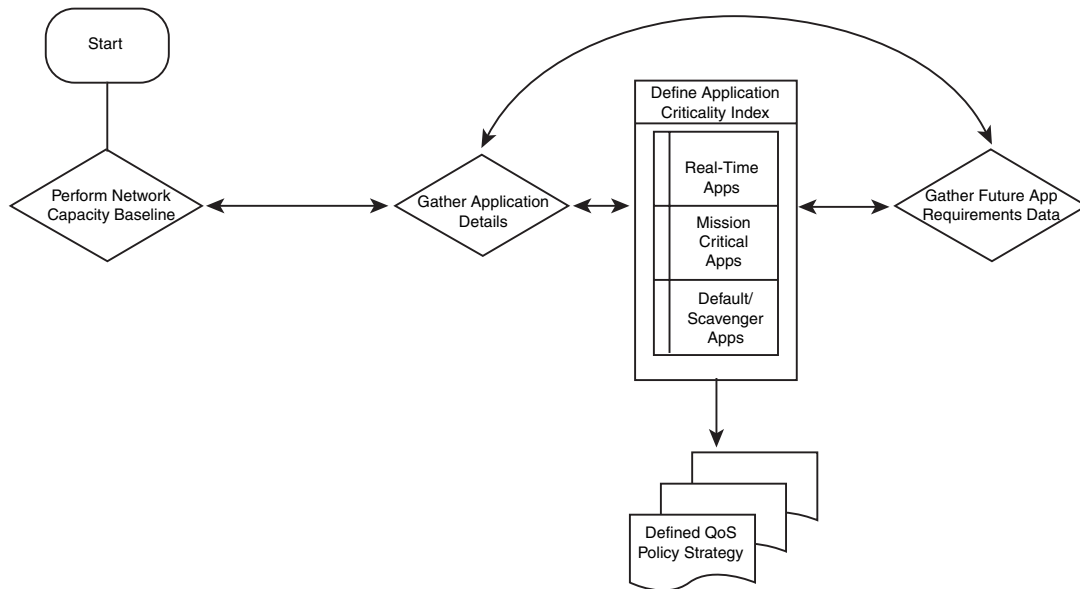
Applications, such as VoIP, videoconferencing over IP, and certain business-critical applications or systems processes, can be one such classification. These applications could be categorized and then assigned a specific handling criteria, such as SAP enterprise resource planning (ERP), storage area network (SAN) replications, Citrix applications, CAD/CAM operations, and, of course, those real-time requirements of voice and video.

The treatment of the traffic is, within a framework of QoS policy, done based on its classification. Endpoints that perform this function need to be able to mark their traffic in a specific way to allow the proper handling to be done as traffic traverses the network. When it comes to real-time applications, this requires identifying the bearer and control (signaling) traffic and ensuring that it is placed in the appropriate class. When this occurs, action needs to be taken to understand what happens with the remaining traffic. Thus, a process of identification is required. It is safe to assume that most applications that traverse the network are, for the most part, unclassified or unmarked where no QoS policy is currently in play.

Figure 5-12 takes the simplistic approach of baselining the existing network to determine the makeup of traffic and applications in play. This is useful because it helps serve as the checkpoint from where to start the classification work. It also requires that an audit be done of the existing network infrastructure to assess its ability to support the applicable policies.

Areas such as device capability and Cisco IOS version must be brought to a consistent level to ensure that there is a base level of capability, stability, and ability to execute the policy.

Figure 5-12 *Identifying the Policy Process*



There are many different mechanisms to start to derive the application communications flow, the most useful of which is Cisco NetFlow Accounting. NetFlow provides valuable information about who is using the network, what applications are being used, when the network is used, and where traffic is going on the network. Essentially it is a way to answer the questions of who, where, when, and what. The mechanism extracts the following:

- Source IP address
- Destination IP address
- Source port
- Destination port
- Layer 3 protocol type
- Type of service byte (DSCP)
- Input logical interface (ifIndex)

Using this capability helps you build a scalable mechanism by which to classify application characteristics and start the work of placing these within a classification framework.

QoS Requirements for Voice, Video, and Data

The next step in the process is to identify the detailed requirements for the classes you would create as part of the QoS framework. These characteristics are well defined, and the process involved in their identification is the critical component to ensuring a consistent approach.

To achieve such values, enterprises and service providers must cooperate and be consistent in classifying, provisioning, and integrating their respective QoS solutions.

Additionally, the service provider's network must support a bare minimum of three classes at all interfaces. This is to ensure that the enterprise need not modify or rework its policy to map to the SP policy needlessly or carry the classes within its own classes transparently. These classes must include a real-time class using LLQ, a high-priority data class with CBWFQ's "minimum bandwidth," and a best-effort class.

Four or five classes are preferred, with the minimum requirements for an LLQ class and all but one of the remaining classes supporting minimum bandwidth.

QoS requirements and high-level recommendations for voice, video, and data are outlined in the following sections.

QoS Requirements for Voice

Voice calls, either one-to-one or on a conference connection capability, require the following:

- ≤ 150 ms of one-way latency from mouth to ear (per the ITU G.114 standard)
- ≤ 30 ms jitter
- ≤ 1 percent packet loss
- 17 to 106 kbps of guaranteed priority bandwidth per call (depending on the sampling rate, codec, and Layer 2 overhead)
- 150 bps (plus Layer 2 overhead) per phone of guaranteed bandwidth for voice control traffic

The choice of codec has impacts in many areas. The most important is the capacity planning on the network, because the bandwidth consumed in different codecs varies.

When exploring the details of these needs in their work on tight IP SLA, John Evans and Clarence Filisfilis wrote that G.114 states that 150 ms of end-to-end one-way delay does not cause a perceivable degradation in voice quality for most use of telephony.

These targets typically include a U.S. coast-to-coast call (equivalent to a Pan-European call) of 6000 km at a propagation speed of 200,000 km/s—thus, 30 ms.

Some carriers try to push to the 100-ms target (excellent: 70 ms without propagation).

A usual target is 150 ms (good: 120 ms without propagation).

Enterprise VoIP networks tend to have a looser target—250 ms (a decent limit: 220 ms without propagation).

It is also recommended that you look at the consumption of Layer 2 overhead; an accurate method for provisioning VoIP is to include the Layer 2 overhead. Layer 2 overhead includes preambles, headers, flags, cyclic redundancy checks (CRCs), and ATM cell padding. When Layer 2 overhead is included in the bandwidth calculations, the VoIP call bandwidth needs translate to the requirements shown in Table 5-3.

Table 5-3 *VoIP Bandwidth Reference Table*

Codec	Sampling Rate	Voice Payload in Bytes	Packets per Second	Bandwidth per Conversation
G.711	20 ms	160	50	80 kbps
G.711	30 ms	240	33	74 kbps
G.729A	20 ms	20	50	24 kbps
G.729A	30 ms	30	33	19 kbps

A more accurate method for the provisioning is to include the Layer 2 overhead in the bandwidth calculations, as shown in Table 5-4.

Table 5-4 *VoIP Bandwidth Needs with Layer 2 Overhead*

Codec	801.Q Ethernet + 32 Layer 2 Bytes	MLP + 13 Layer 2 Bytes	Frame Relay + 8 Layer 2 Bytes	ATM + Variable Layer 2 Bytes (Cell Padding)
G.711 at 50 pps	93 kbps	86 kbps	84 kbps	104 kbps
G.711 at 33 pps	83 kbps	78 kbps	77 kbps	84 kbps
G.711 at 50 pps	37 kbps	30 kbps	28 kbps	43 kbps
G.711 at 33 pps	27 kbps	22 kbps	21 kbps	28 kbps

Sample Calculation

The following calculations are used to determine the inputs to the planning of voice call consumption:

total packet size = (L2 header: MP or FRF.12 or Ethernet) + (IP/UDP/RTP header) + (voice payload size)

pps = (codec bit rate) / (voice payload size)

bandwidth = total packet size * pps

For example, the required bandwidth for a G.729 call (8-kbps codec bit rate) with cRTP, MP, and the default 20 bytes of voice payload is as follows:

$$\text{total packet size (bytes)} = (\text{MP header of 6 bytes}) + (\text{compressed IP/UDP/RTP header of 2 bytes}) + (\text{voice payload of 20 bytes}) = 28 \text{ bytes}$$
$$\text{total packet size (bits)} = (28 \text{ bytes}) * 8 \text{ bits per byte} = 224 \text{ bits}$$
$$\text{pps} = (8\text{-kbps codec bit rate}) / (160 \text{ bits}) = 50 \text{ pps}$$

NOTE 160 bits = 20 bytes (default voice payload) * 8 bits per byte

$$\text{bandwidth per call} = \text{voice packet size (224 bits)} * 50 \text{ pps} = 11.2 \text{ kbps}$$

QoS Requirements for Video

The requirements for streaming video, such as IP multicast, executive broadcasts, and real-time training activities, are as follows:

- Four to 5 seconds of latency allowable (depending on the video application's buffering capabilities). No significant jitter requirements.
- Two percent packet loss permissible. Bandwidth required depends on the encoding and the rate of the video stream.
- Video content distribution such as video on demand being replicated to distributed content engines.
- Delay- and jitter-insensitive.
- Large file transfers (traffic patterns similar to FTP sessions).
- Restrict to distribution to less-busy times of day.
- Provision as "less-than-best-effort" data.

The requirements for videoconferencing can be applied as either a one-to-one capability or a multipoint conference.

- ≤ 150 ms of one-way latency from mouth to ear (per the ITU G.114 standard).
- ≤ 30 ms jitter.
- ≤ 1 percent packet loss.
- Minimum bandwidth guarantee is videoconferencing session + 20 percent. For example, a 384-kbps videoconferencing session requires 460 kbps guaranteed priority bandwidth.

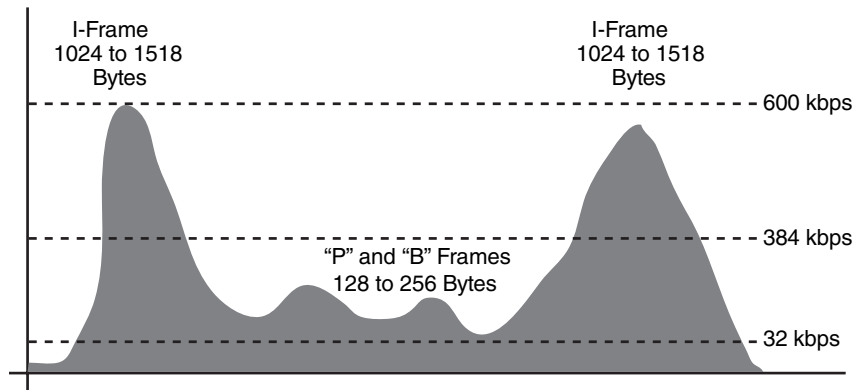
I-Frames are full-frame samples, whereas P and B frames are differential (or delta) frames. Videoconferencing shares the same latency, jitter, and loss requirements as voice but has radically burstier and heavier traffic patterns.

A 384-kbps stream can take up to 600 kbps at points rather than provisioning the stream + 60 percent (to accommodate the occasional 600-kbps burst).

The video stream includes an additional 20 percent of bandwidth with a burst allowance in the LLQ of 30,000 bytes per 384-kbps stream, as shown in Figure 5-13:

- Provision LLQ to stream + 20 percent.
For example, 384-kbps stream → 460-kbps LLQ
- Additionally, extend the LLQ burst to capture I frames without requiring additional LLQ bandwidth reservation.

Figure 5-13 Video Stream Sequence



QoS Requirements for Data

Following the earlier discussion about the application of traffic identification, there are a few key points to remember about classifying data traffic:

- Profile applications to get a basic understanding of their network requirements.
- Perform capacity planning to ensure an adequate bandwidth baseline.
- Use no more than four separate data traffic classes:
 - Transactional data (mission-critical)—ERP, transactional, and high-priority internal applications

- Bulk data (guaranteed-bandwidth)—Streaming video, messaging, and intranet
- Best-effort (the default class)—Internet browsing, e-mail, and unclassified applications
- Scavenger (less-than-best-effort)—FTP, backups, and noncritical applications
- Minimize the number of applications assigned to the transactional and bulk data classes (three or fewer are recommended).
- Use proactive provisioning policies before reactive policing policies.

These requirements for the data classes are guidelines. They need to take into account many different factors, including the service provider. They must be able to support the number of classes required by the enterprise. As such, they may affect the decision process in the policy's creation.

Governance plays a key role in identifying and classifying applications. By building a governance checkpoint in the development of new applications, a lot of cycles can be reduced in determining an application's bandwidth needs and its impact on network requirements. The process can pinpoint whether the new application will lead to any network upgrades as well as determine a baseline for the application, which can lead to predictive planning on an incremental yearly basis of adding the application to the network. It allows for better planning, thereby removing the challenge that can be created because of bottlenecks or scaling issues, and the key tenant of running the network as a fairly used business asset.

A critical factor in the service provider delivery is the SLA, affecting the ability to support delay-sensitive classes, as seen in voice and video requirements. In some cases, other factors preclude the service provider's ability to deliver an SLA against these requirements, such as the geographic location of sites and the interconnections over the service provider network.

For example, suppose you have a site in a geographically remote location that has a large propagation delay imposed on it because of its location. It may not be possible to meet the requirements for delivery of real-time services to its location. In such cases, there is a trade-off between what is possible for the majority of the corporate sites and that of the remote geographies and their interconnection capability to the rest of the network.

The LAN Edge: L2 Configurations

Framing this in the context of a sample configuration, the switch device at the edge of the network has the following policy applied:

Depending on the switch model, it may be necessary to first activate QoS using this command:

```
switch(config)#mls qos
```

This command is required on both the Catalyst 3550 and the Catalyst 6500. The Catalyst 2950 has QoS enabled by default.

The trust is configured on the switch port using this command:

```
switch(config-if)#mls qos trust dscp
```

Any ISL or 802.1Q/p frames that enter the switch port have their CoS passed (untouched) through the switch. If an untagged frame arrives at the switch port, the switch assigns a default CoS to the frame before forwarding it. By default, untagged frames are assigned a CoS of 0. This can be changed using this interface configuration command:

```
switch(config-if)#mls qos cos default-cos
```

where *default-cos* is a number between 0 and 7.

The syntax to configure QoS trust switch-wide for IP phone endpoints is all that is required in typical installations:

```
Switch(config-if)#mls qos trust device ciscoipphone
```

Here's the legacy syntax that was required on a per-VLAN basis:

```
Switch(config-if)#switchport voice vlan {vlan-id | dot1p | none | untagged}
```

To instruct the Cisco IP Phone to forward all voice traffic through a specified VLAN, use this command:

```
Switch(config-if)#switchport voice vlan vlan-id
```

By default, the Cisco IP Phone forwards the voice traffic with an 802.1Q priority of 5. Valid VLAN IDs are from 1 to 4096.

An alternative to specifying a particular voice VLAN on the switch is to instruct the switch port to use 802.1P priority tagging for voice traffic and to use the default native VLAN (VLAN 0) to carry all traffic. By default, if enabled, the Cisco IP Phone forwards the voice traffic with an 802.1P priority of 5.

```
Switch(config-if)#switchport voice vlan dot1p
```

In some cases, it may be desirable—indeed, highly recommended—not to trust edge CoS for nonvoice/video endpoints and not to trust any CoS value that may be present in frames sourced from an edge device. For example, an office PC used for general applications, such as web browsing, e-mail, and file and print services, may not require special QoS treatment. Allowing it to request higher levels of QoS may adversely affect applications such as voice and video, which require guarantees of bandwidth and latency.

NOTE

This may not hold true for data center or server-based systems, which need to be given individual consideration based on the prioritization needs of any application they serve.

For this reason, it is possible to use the **override** parameter to tell the switch to ignore any existing CoS value that may be in the frame and apply the default value. This effectively disables any trust configuration that may have previously been applied to the port.

The CoS value assigned by the switch can be changed on a port-by-port basis using this interface configuration command:

```
Switch(config-if)#mls qos cos override
```

After this command is applied, the switch rewrites the CoS value for all incoming frames to the configured default value, regardless of any existing CoS value.

Other platforms, such as those that employ CatOS, vary. You should always verify such a reference to the correct procedure by reviewing the relevant documentation at <http://www.cisco.com>. For example, the following is an overview of configuring prioritization, with a 6500 platform running CatOS between Cisco CallManager and IP phones and gateways using TCP ports 2000 to 2002. The sample commands classify all Skinny Protocol traffic from IP phones and gateways (VLAN 110) and Cisco CallManager (4/2) as DSCP AF31, which is backward-compatible with IPP 3.

With older implementations, several steps need to be performed (see Example 5-1):

- Step 1** Enable switch-wide QoS.
- Step 2** Create an ACL (ACL_IP-PHONES), marking all Skinny Client and Gateway Protocol traffic from the IP phones and from Skinny Protocol gateways with a DSCP value of AF31.
- Step 3** Add to the ACL_IP-PHONE access list, trusting all DSCP markings from the IP phone so that the IP Prec = 5 RTP traffic is not rewritten.
- Step 4** Create an ACL (ACL_VOIP_CONTROL), marking all Skinny Client and Gateway Protocol traffic from Cisco CallManager with a DSCP value of AF31.
- Step 5** Accept incoming Layer 2 CoS classification.
- Step 6** Inform the port that all QoS associated with the port will be done on a VLAN basis to simplify configuration.
- Step 7** Instruct the IP phone to rewrite CoS from the PC to CoS=0 within the IP phone Ethernet ASIC.
- Step 8** Inform Cisco CallManager port (4/2) that all QoS associated with the port will be done on a port basis.
- Step 9** Write the ACL to hardware.
- Step 10** Map the ACL_IP-PHONE ACL to the auxiliary VLAN.

Step 11 Map the ACL_VOIP_CONTROL ACL to the Cisco CallManager port.

Example 5-1 Setup of a Catalyst Edge Switch (L2 Only)

```

cat6k-access> (enable) set qos enable
cat6k-access> (enable) set qos acl ip ACL_IP-PHONES dscp 26 tcp any any range 2000
2002
cat6k-access> (enable) set qos acl ip ACL_IP-PHONES trust-cos ip any any
cat6k-access> (enable) set qos acl ip ACL_VOIP_CONTROL dscp 26 tcp any any range
2000 2002
cat6k-access> (enable) set port qos 5/1-48 trust trust-cos
cat6k-access> (enable) set port qos 5/1-48 vlan-based
cat6k-access> (enable) set port qos 5/1-48 trust-ext untrusted
cat6k-access> (enable) set port qos 4/2 port-based
cat6k-access> (enable) commit qos acl all
cat6k-access> (enable) set qos acl map ACL_IP-PHONES 110
cat6k-access> (enable) set qos acl map ACL_VOIP_CONTROL 4/2

```

Classifying Voice on the WAN Edge

A basic configuration of a WAN edge (CE) router is shown in Example 5-2 and defined further in this section. This applies the principles in base configuration terms, as discussed previously.

In this case, you apply a simple LLQ and CBWFQ policy to the router to support voice traffic. Voice traffic needs to be assigned to the LLQ, and voice-control traffic needs a minimum bandwidth guarantee.

Example 5-2 Matching Voice and Voice Control

```

ip cef
!
class-map match-all VOICE
  match ip dscp ef
!
class-map match-all VOICE-CONTROL
  match ip dscp cs3
  match ip dscp af31
!
!
policy-map WAN-EDGE
  class VOICE
    priority percent 33
  class VOICE-CONTROL
    bandwidth percent 5
  class class-default
    fair-queue
!
interface Serial0/0
  description WAN Link to CCT ID : 1234 :: SP-PE-1
  bandwidth 2048
  ip address 10.1.1.1 255.255.255.252
  service-policy output WAN-EDGE

```

The class map applied to VOICE and VOICE-CONTROL provides an example of matching against DSCP markings, which in this case are for the voice-bearer traffic and the voice signaling and control traffic. In this basic policy, you assume that no marking of traffic needs to happen directly on this box. Rather, it deals only with the outbound classification of the traffic.

In the policy map itself, in effect you assign three elements—those of prioritization for voice and specified as a percentage of the link bandwidth. This is an example of applying LLQ. Voice control specifies the allocation of 5 percent of link bandwidth for voice control, with the final classification being performed as class default, which uses the fair-queuing algorithm on the remaining traffic. In this case, it is assumed to be the default class.

The key to activating this policy is to attach the service-policy output to an interface—in this case, the output interface facing the SP PE. Thus, applying the classification actions as defined in the policy map WAN-EDGE to the output of the WAN serial interface provides this activation.

Classifying Video on the WAN Edge

In Example 5-3, videoconferencing traffic is assigned LLQ, and all nonvideo traffic is assigned to a default queue for WFQ.

Example 5-3 *Matching Video*

```
ip cef
!
class-map match-all VIDEO
  match ip dscp af41
!
!
policy-map WAN-EDGE
  class VIDEO-CONF
    priority 460
  class class-default
    fair-queue
!
interface Serial0/0
  description WAN Link to CCT ID : 1234 :: SP-PE-1
  bandwidth 2048
  ip address 10.1.1.1 255.255.255.252
  service-policy output WAN-EDGE
```

On the WAN edge, videoconferencing traffic should be assigned to an LLQ. The video stream minimum bandwidth guarantee should be the size of the stream plus 20 percent.

As before, this policy doesn't take effect until it is bound to an interface with a service-policy statement.

Classifying Data on the WAN Edge

Most enterprises have many applications that can be considered mission-critical (transactional). However, if too many applications are classified as mission-critical, they will contend among themselves for bandwidth, with the result of dampening QoS effectiveness. Taken to the extreme, a regular FIFO link (no QoS) is scheduled in the exact same manner as a link where every application is provisioned as mission-critical. Therefore, it is recommended that you classify no more than three applications as mission-critical (transactional).

These applications should be marked with different AF drop-preference values to distinguish them from each other. Such distinctions provide more granular visibility in managing and monitoring application traffic and aid in provisioning for future requirements. Similar arguments are made for having no more than three applications in a guaranteed bandwidth (bulk data) class of applications and, likewise, marking these applications with different AF drop-preference values.

Default traffic is automatically marked as best-effort (DSCP 0). However, noncritical bandwidth-intensive traffic could (optionally) be marked as different so that adverse policies could be applied to control such traffic. These types of traffic can be described as “less-than-best-effort” or “scavenger” traffic.

It is imperative that DSCP classification be performed on all packets before they arrive at the WAN edges. In this manner, queuing and congestion avoidance can be performed at the WAN edge based strictly on DSCP markings.

NOTE

It’s important to keep in mind that the default class map **match** setting is **match-all**. Therefore, when you attempt to classify mutually exclusive traffic flows (such as differing DSCP values), it is important to explicitly use the **match-any** qualifier when defining the class map. Another example could be through the use of multiple DSCPs on a single command line:

```
match ip dscp af11 af12 af13
```

The advantage of using multiple lines is that this triggers a separate counter in the class-based QoS Management Information Base (MIB) for each DSCP. (If they are matched all on the same line, only a single counter is triggered for all DSCPs.)

The Eight-Class Model introduces a dual-LLQ design: one for voice and another for interactive video. As pointed out earlier in this chapter, the LLQ has an implicit policer that allows for time-division multiplexing of the single priority queue. This implicit policer abstracts the fact that there is essentially a single LLQ within the algorithm and thus allows for the “provisioning” of multiple LLQs.

Interactive video (or IP videoconferencing, also called IP/VC) is recommended to be marked AF41 (which can be marked down to AF42 in the case of dual-rate policing at the campus access edge). It is recommended that you overprovision the LLQ by 20 percent of the IP/VC rate. This takes into account IP/UDP/RTP headers as well as Layer 2 overhead. Additionally, Cisco IOS Software automatically includes a 200-ms burst parameter (defined in bytes) as part of the priority command. On dual-T1 links, this has proven sufficient for protecting a single 384-kbps IP/VC stream. On higher-speed links (such as triple T1s), the default burst parameter has shown to be insufficient for protecting multiple IP/VC streams.

However, multiple-stream IP/VC quality tested well with the burst set to 30,000 bytes (for example, priority 920 30000). Our testing did not arrive at a clean formula for predicting the required size of the burst parameters as IP/VC streams continually were added. However, given the variable packet sizes and rates of these interactive-video streams, this is not surprising. The main point is that the default LLQ burst parameter might require tuning as multiple IP/VC streams are added (which likely will be a trial-and-error process).

Optionally, DSCP-based WRED can be enabled on the Interactive-Video class, but testing has shown negligible performance difference in doing so (because, as has been noted, WRED is more effective on TCP-based flows than UDP-based flows, such as interactive video). In these designs, WRED is not enabled on classes such as Call-Signaling, IP Routing, and Network-Management because WRED would take effect only if such classes were filling their queues nearly to their limits. Such conditions would indicate a provisioning problem that would be better addressed by increasing the class's minimum bandwidth allocation than by enabling WRED.

Additionally, the Eight-Class Model subdivides the preferential data class to separate control plane traffic (IP routing and network-management applications) from business-critical data traffic. IGP packets (such as RIP, EIGRP, OSPF, and IS-IS) are protected through the PAK_priority mechanism within the router. However, EGP protocols, such as BGP, do not get PAK_priority treatment and might need explicit bandwidth guarantees to ensure that peering sessions do not reset during periods of congestion. Additionally, administrators might want to protect network-management access to devices during periods of congestion.

The other class added to this model is for bulk traffic (the Bulk Data class), which is spun off of the Critical Data class. Because TCP continually increases its window sizes, which is especially noticeable in long sessions (such as large file transfers), constraining bulk data to its own class alleviates other data classes from being dominated by such large file transfers. Bulk data is identified by DSCP AF11 (or AF12 in the case of dual-rate policing at the campus access edges). DSCP-based WRED can be enabled on the Bulk Data class

(and also on the Critical Data class). Example 5-4 shows the implementation of the Eight-Class Model for the WAN edge.

Example 5-4 *Eight-Class Model*

```

!
class-map match-all Voice
  match ip dscp ef ! IP Phones mark Voice to EF
class-map match-all Interactive Video
  match ip dscp af41 af42 ! Recommended markings for IP/VC
class-map match-any Call Signaling
  match ip dscp cs3 ! Future Call-Signaling marking
  match ip dscp af31 ! Current Call-Signaling marking
class-map match-any Network Control
  match ip dscp cs6 ! Routers mark Routing traffic to CS6
  match ip dscp cs2 ! Recommended marking for Network Management
class-map match-all Critical Data
  match ip dscp af21 af22 ! Recommended markings for Transactional-Data
class-map match-all Bulk Data
  match ip dscp af11 af12 ! Recommended markings for Bulk-Data
class-map match-all Scavenger
  match ip dscp cs1 ! Scavenger marking
!
policy-map WAN-EDGE
  class Voice
    priority percent 18 ! Voice gets 552 kbps of LLQ
  class Interactive Video
    priority percent 15 ! 384 kbps IP/VC needs 460 kbps of LLQ
  class Call Signaling
    bandwidth percent 5 ! BW guarantee for Call-Signaling
  class Network Control
    bandwidth percent 5 ! Routing and Network Management get min 5% Bandwidth
  class Critical Data
    bandwidth percent 27 ! Critical Data gets min 27% BW
    random-detect dscp-based ! Enables DSCP-WRED for Critical-Data class
  class Bulk Data
    bandwidth percent 4 ! Bulk Data gets min 4% BW guarantee
    random-detect dscp-based ! Enables DSCP-WRED for Bulk-Data class
  class Scavenger
    bandwidth percent 1 ! Scavenger class is throttled to 1% of bandwidth
  class class-default
    bandwidth percent 25 ! Fair-queuing is sacrificed for Bandwidth guarantee
    random-detect ! Enables WRED on class-default
!

```

This design is more efficient than strict policing, because these bandwidth-intensive noncritical applications can use additional bandwidth if it is available while ensuring protection for critical transactional data classes and real-time voice and video applications.

The development of this model is attributed to Tim Szigeti from the Cisco Enterprise Solutions Engineering Group.

For more information on this approach, go to http://www.cisco.com/application/pdf/en/us/guest/netsol/ns432/c649/ccmigration_09186a008049b062.pdf.

Case Study: QoS in the Acme, Inc. Network

Acme, Inc. currently uses five different classes of service on the Acme network. For simplicity's sake, and because of the 3-bit resolution of some QoS technologies, such as 802.1p LAN CoS, MPLS EXP, and IP precedence, Acme IT limits its QoS to classes that can be identified in 3 bits. Traffic is classified at the WAN edge by matching IP precedence values—only the first 3 bits of the ToS byte, which cover 7 DSCP values each.

The policy template shown in Table 5-5 deals with policy marking of packets in Acme. It shows QoS policy for voice, video, and data classifications and their breakout assignments based on circuit speed, allocation per class as a percentage, and the classifications assigned.

Table 5-5 *Sample Policy Template Breakdown*

Policy Number	30	25	20	15	10	5
	Bandwidth at Interface Line Rate (in kbps)					
Class	622000	155000	45000 to 34000	34000 to 2048	2048 to 1024	1024 to 256
Management	2%	2%	2%	2%	2%	6%
Voice EF	10%	33%	33%	33%	33%	33%
Video AF4	5%	5%	8%	7%	10%	5%
Signaling AF3	5%	5%	5%	5%	5%	10%
Default BE	58%	45%	50%	43%	40%	36%
Scavenger CS1	20%	10%	10%	10%	10%	10%
Total	100%	100%	100%	100%	100%	100%

IP precedence values 6 and 7 are also used for network control traffic (routing protocols). Most of this traffic is link-local (CE-PE only), allowing an individual class of traffic to be set up for this traffic on a WAN port with minimum bandwidth. On the CE side of the CE-to-PE links, it is recommended that a separate class be used for management traffic. On the

PE side of the CE-to-PE link, this tends to vary with each provider. This traffic must, at a minimum, be mapped to a high-priority data class of service in the service provider cloud.

QoS for Low-Speed Links: 64 kbps to 1024 kbps

LFI allows large packets on a serial link to be divided into using MLP or Frame Relay encapsulation with FRF.12 fragmentation.

To determine the LFI fragment size, you must consider the packet flow through the router. Following the link fragmentation process and LLQ/CBWFQ's fragment ordering, fragments are placed on a transmit ring buffer or TX ring. The TX ring then queues packets onto the physical interface in a FIFO fashion.

Example 5-5 shows examples of applications using MLP and FRF.12.

Example 5-5 LFI on MLP

```

interface Multilink1                               Multilink bundle interface
bandwidthd XXX                                    Enter Link bandwidth in kbps
ip address 10.52.255.1 255.255.255.252
no ip redirects
no ip proxy-arp
ip authentication mode eigrp 109 md5
ip authentication key-chain eigrp 100 apple_key
max-reserved-bandwidth 100
service-policy output WAN-EDGE                    Apply service policy outbound
ppp multilink                                       Configure Multilink
ppp multilink fragment-delay 10                   Set the max packet delay in ms
                                                    (determines fragment size)
ppp multilink interleave                           Enable LFI
multilink-group 1                                  Apply template to multilink group #1
!
!
interface Serial S:P
description Multilink PPP group member
bandwidthd XXX                                    Configure bandwidth equal to full line rate
no ip address
no ip redirects
no ip proxy-arp
encapsulation ppp
fair-queue
ppp multilink                                       enable multilink on interface
multilink-group 1                                  Assign interface to multilink group 1
!
interface Serial S:P                               Next I/F when MLP Bundling is required
description Multilink PPP group member
bandwidthd XXX

```

Slow-Speed (768-kbps) Leased-Line Recommendation: Use MLP LFI and cRTP

For slow-speed leased lines, LFI is required to minimize serialization delay. Therefore, MLP is the only encapsulation option on slow-speed leased lines because MLP LFI is the only mechanism available for fragmentation and interleaving on such links. Optionally, cRTP can be enabled either as part of the modular QoS command-line interface (MQC) policy map or under the multilink interface (using the **ip rtp header-compression** command). Ensure that MLP LFI and cRTP, if enabled, are configured on both ends of the point-to-point link, as shown in Example 5-6.

Example 5-6 *Slow-Speed (768-kbps) Leased-Line QoS Design Example*

```

!
policy-map WAN-EDGE
class Voice
priority percent 33 ! Maximum recommended LLQ value
compress header ip rtp ! Enables Class-Based cRTP
class Call Signaling
bandwidth percent 5 ! BW guarantee for Call-Signaling
!
interface Multilink1
description 768 kbps Leased-Line to RBR-3745-Left
ip address 10.1.112.1 255.255.255.252
service-policy output WAN-EDGE ! Attaches the MQC policy to Mu1
ppp multilink
ppp multilink fragment delay 10 ! Limits serialization delay to 10 ms
ppp multilink interleave ! Enables interleaving of Voice with Data
ppp multilink group 1
!
...
!
interface Serial1/0
bandwidth 786
no ip address
encapsulation ppp
ppp multilink
ppp multilink group 1 ! Includes interface Ser1/0 into Mu1 group
!

```

These examples cover the application of the WAN-EDGE service policy discussed in Example 5-4. For more examples of configuring the WAN edge, refer to the Cisco QoS Solution Reference Network Design, http://www.cisco.com/application/pdf/en/us/guest/netso/ns432/c649/ccmigration_09186a008049b062.pdf.

QoS Reporting

The service provider should offer a minimum level of QoS reporting to provide you with statistics on the QoS you are receiving. This reporting should include the following metrics:

- Observed traffic ratio for classes
- Number of packet drops per class as a number and percentage
- Average offered load per class

The network administrator may specify additional parameters as necessary to ensure the network's maintenance and operation.

You can use the Cisco IP SLA Response Time Reporter (RTR) feature in Cisco IOS to measure the response time between IP devices. A source router configured with IP SLA configured can measure the response time to a destination IP device that can be a router or an IP device. The response time can be measured between the source and the destination or for each hop along the path. Simple Network Management Protocol (SNMP) traps can be configured to alert management consoles if the response time exceeds the predefined thresholds.

Key areas that can be measured with IP SLA include

- Interpacket delay variance (jitter) of VoIP traffic
- Response time between endpoints for a specific QoS
- IP ToS bits
- Packet loss using IP SLA generated packets

You can configure the IP SLA feature on routers using the Cisco Internetwork Performance Monitor (IPM) application. The IP SLA/RTR is imbedded in many but not all feature sets of the Cisco IOS software. A release of the Cisco IOS software that supports IP SLA/RTR must be installed on the device that IPM uses to collect performance statistics.

Chapter 8, "Network Management, SLA, and Support Requirements," contains more information on reporting, monitoring, and managing QoS.

Summary

This chapter assessed the need for QoS and the implications for the enterprise when you use an MPLS VPN-based service. The conclusion is that the need for QoS is real in the LAN, WAN, and SP transit network in support of the move toward a converged network. You must understand the mechanisms that are available to solve the application of QoS in the network, how to take a rational and staged approach toward developing a scalable policy, and how to execute that policy.

QoS is not a silver bullet that creates bandwidth. It needs to be used as part of a well-defined capacity planning and overall application governance process that identifies current and future state evolutions. Ensuring that the network can support the enterprise developments in a productive and manageable form should be the aim of any policy development so that it remains manageable over time.

The common theme is planning, designing, implementing, and operating. Build on these themes to ensure that there is consistency in the QoS strategy from an end-to-end architecture perspective. As this technology evolves, there will be additional areas to explore, especially call admission control and resource reservation. Indeed, the recent implementation of “MLPP for Voice and Video in the Internet Protocol Suite,” currently in the IETF Draft stage by F. Baker, et al., shows promise as to developing the right answers.

To more fully explore the vast topic of QoS, it is highly recommended that you read the Cisco Press book *End-to-End QoS Network Design: Quality of Service in LANs, WANs, and VPNs* by Tim Szigeti and Christina Hattingh.

References

RFC 2210, *The Use of RSVP with IETF Integrated Services*

<http://www.faqs.org/rfcs/rfc2210.html>

RFC 2309, *Recommendations on Queue Management and Congestion Avoidance in the Internet*

<http://www.faqs.org/rfcs/rfc2309.html>

RFC 2474, *Definition of the Differentiated Services Field (DS Field) in the IPv4 and IPv6 Headers*

<http://www.faqs.org/rfcs/rfc2474.html>

RFC 2475, *An Architecture for Differentiated Service*

<http://www.faqs.org/rfcs/rfc2475.html>

RFC 2547, *BGP/MPLS VPNs*

<http://www.faqs.org/rfcs/rfc2547.html>

RFC 3168, *The Addition of Explicit Congestion Notification (ECN) to IP*

<http://www.faqs.org/rfcs/rfc3168.html>

RFC 3270, *Multi-Protocol Label Switching (MPLS) Support of Differentiated Services*

<http://www.faqs.org/rfcs/rfc3270.html>

QoS Overview

http://www.cisco.com/univercd/cc/td/doc/cisintwk/ito_doc/qos.htm

Cisco IP Telephony Solution Reference Network Design Guide

http://www.cisco.com/warp/public/779/largeent/netpro/avvid/iptel_register.html

Cisco Solutions Reference Network Design (SRND) Guide for QoS

http://www.cisco.com/application/pdf/en/us/guest/netsol/ns17/c649/ccmigration_09186a00800d67ed.pdf

IP Videoconferencing Solution Reference Network Design Guide

http://www.cisco.com/warp/public/779/largeent/netpro/avvid/ipvc_register.html

Low Latency Queuing with Priority Percentage Support

<http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122newft/122t/122t2/ftllqpc.htm>

Class-Based Marking

<http://www.cisco.com/univercd/cc/td/doc/product/software/ios121/121newft/121t/121t5/cbpmark2.htm>

Configuring Frame Relay Traffic Shaping

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fwan_c/wcffrely.htm#xtocid27

Configuring Distributed Traffic Shaping

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcprt4/qcfdts.htm

Configuring ATM

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fwan_c/wcfatm.htm

IP Header Compression Enhancement—PPPoATM and PPPoFR Support

<http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122relnt/xprn122t/122tnewf.htm#xtocid274>

Configuring Frame Relay-ATM Interworking

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fwan_c/wcffratm.htm

Classification in the Cisco IOS Quality of Service Solutions Configuration Guide, Release 12.3

http://www.cisco.com/univercd/cc/td/doc/product/software/ios123/123cgcr/qos_vcg.htm

Service Provider QoS Design Guide

http://www.cisco.com/en/US/netsol/ns341/ns396/ns172/ns103/networking_solutions_white_paper09186a00801b1c5a.shtml

Network-Based Application Recognition

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcprt1/qcfclass.htm#xtocid24
<http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122newft/122t/122t8/dtnbarad.htm>

Congestion Management in the Cisco IOS Quality of Service Solutions Configuration Guide, Release 12.2

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcprt2/qcfconmg.htm

Congestion Avoidance in the Cisco IOS Quality of Service Solutions Configuration Guide, Release 12.3

http://www.cisco.com/univercd/cc/td/doc/product/software/ios123/123cgcr/qos_vcg.htm#1000448

Policing and Shaping in the Cisco IOS Quality of Service Solutions Configuration Guide, Release 12.2

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcprt4/index.htm

Link Efficiency Mechanisms in the Cisco IOS Quality of Service Solutions Configuration Guide, Release 12.2

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcprt6/index.htm

Configuring Compressed Real-Time Protocol

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcprt6/qcfcrt.htm

Modular Quality of Service Command-Line Interface in the Cisco IOS Quality of Service Solutions Configuration Guide, Release 12.2

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcprt8/index.htm

Traffic Policy as a QoS Policy (Hierarchical Traffic Policies) Example

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcprt8/qcfmcli2.htm#xtocid16

IP Telephony QoS Design Guide

http://www.cisco.com/univercd/cc/td/doc/product/voice/ip_tele/avvidqos/

Cisco Class-Based QoS Configuration and Statistics MIB

<ftp://ftp.cisco.com/pub/mibs/v2/CISCO-CLASS-BASED-QOS-MIB.my>

WRED, Explicit Congestion Notification (ECN)

<http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122newft/122t/122t8/ftwrdecn.htm>

NetFlow

http://www.cisco.com/warp/public/732/Tech/nmp/netflow/netflow_learnabout.shtml

Voice Calculations

http://www.cisco.com/warp/public/788/pkt-voice-general/bwidth_consume.html

Video Calculations

<http://www.cisco.com/warp/public/105/video-qos.html>

H.323 Deployments

http://www.cisco.com/warp/public/cc/pd/iosw/ioft/mmcm/tech/h323_wp.htm

Cisco Products Quick Reference Guide: February 2002

<http://www.cisco.com/warp/public/752/qrg/cpqrg2.htm>