



Module 9: QoS Concepts

Enterprise Networking, Security,
and Automation v7.0 (ENSA)





Module Objectives

Module Title: QoS Concepts

Module Objective: Explain how networking devices implement QoS.

Topic Title	Topic Objective
Network Transmission Quality	Explain how network transmission characteristics impact quality.
Traffic Characteristics	Describe minimum network requirements for voice, video, and data traffic.
Queuing Algorithms	Describe the queuing algorithms used by networking devices.
QoS Models	Describe the different QoS models.
QoS Implementation Techniques	Explain how QoS uses mechanisms to ensure transmission quality.

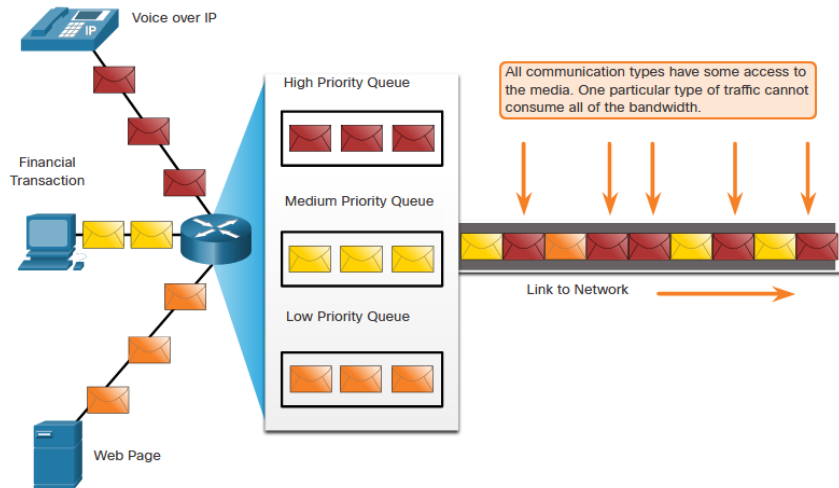


9.1 Network Transmission Quality

Network Transmission Quality

Prioritizing Traffic

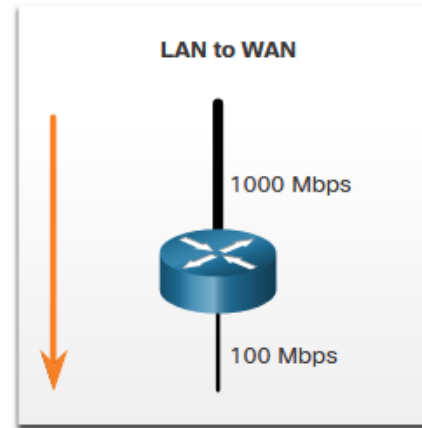
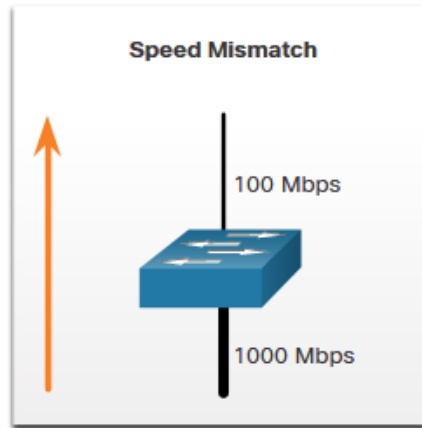
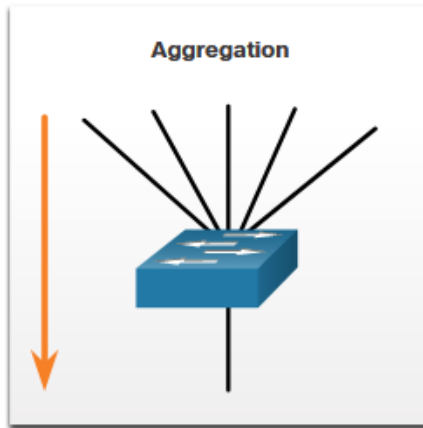
- When traffic volume is greater than what can be transported across the network, devices queue (hold) the packets in memory until resources become available to transmit them.
- Queuing packets causes delay because new packets cannot be transmitted until previous packets have been processed.
- If the number of packets to be queued continues to increase, the memory within the device fills up and packets are dropped.
- One QoS technique that can help with this problem is to classify data into multiple queues, as shown in the figure.



Note: A device implements QoS only when it is experiencing some type of congestion.

Bandwidth, Congestion, Delay, and Jitter

- Network bandwidth is measured in the number of bits that can be transmitted in a single second, or bits per second (bps).
- Network congestion causes delay. An interface experiences congestion when it is presented with more traffic than it can handle. Network congestion points are ideal candidates for QoS mechanisms.
- The typical congestion points are aggregation, speed mismatch, and LAN to WAN.



Bandwidth, Congestion, Delay, and Jitter (Cont.)

Delay or latency refers to the time it takes for a packet to travel from the source to the destination.

- Fixed delay is the amount of time a specific process takes, such as how long it takes to place a bit on the transmission media.
- Variable delay takes an unspecified amount of time and is affected by factors such as how much traffic is being processed.
- Jitter is the variation of delay of received packets.

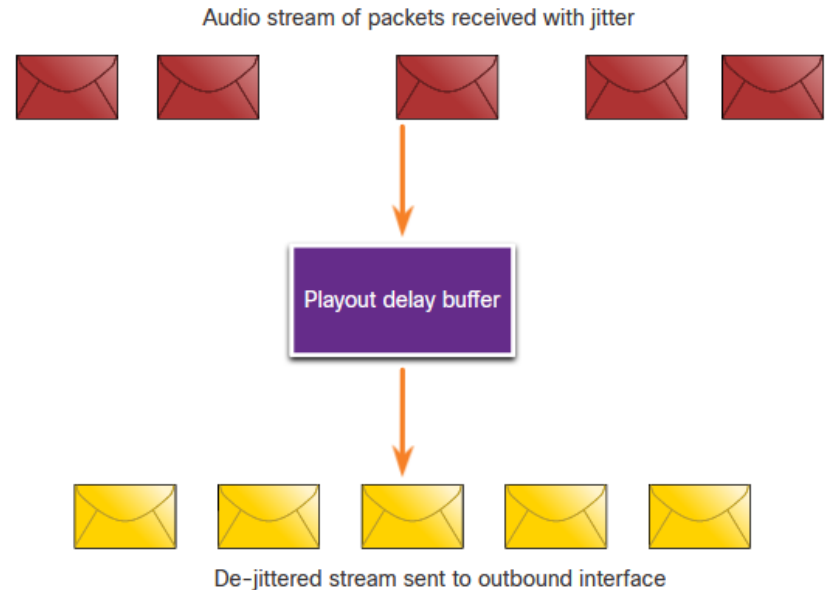
Delay	Description
Code delay	The fixed amount of time it takes to compress data at the source before transmitting to the first internetworking device, usually a switch.
Packetization delay	The fixed time it takes to encapsulate a packet with all the necessary header information.
Queuing delay	The variable amount of time a frame or packet waits to be transmitted on the link.
Serialization delay	The fixed amount of time it takes to transmit a frame onto the wire.
Propagation delay	The variable amount of time it takes for the frame to travel between the source and destination.
De-jitter delay	The fixed amount of time it takes to buffer a flow of packets and then send them out in evenly spaced intervals.

Network Transmission Quality

Packet Loss

Without QoS mechanisms, time-sensitive packets, such as real-time video and voice, are dropped with the same frequency as data that is not time-sensitive.

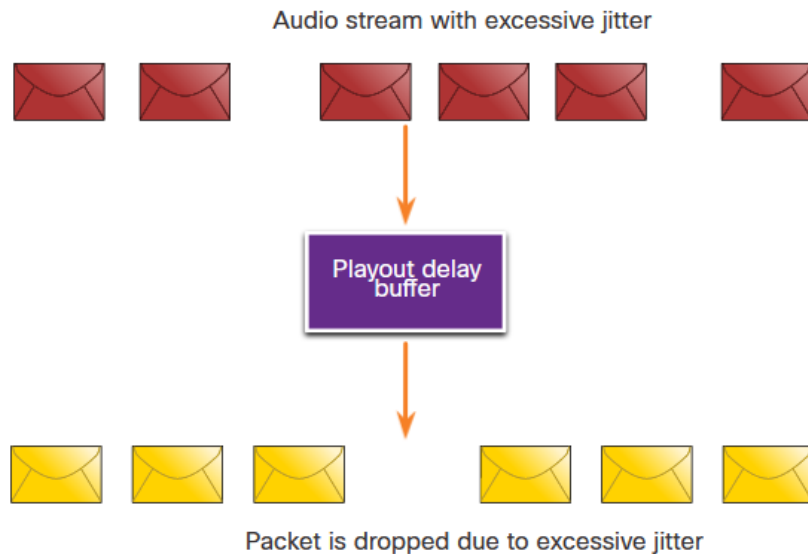
- When a router receives a Real-Time Protocol (RTP) digital audio stream for Voice over IP (VoIP), it compensates for the jitter that is encountered using a playout delay buffer.
- The playout delay buffer buffers these packets and then plays them out in a steady stream.



Network Transmission Quality Packet Loss (Cont.)

If the jitter is so large that it causes packets to be received out of the range of the play out buffer, the out-of-range packets are discarded and dropouts are heard in the audio.

- For losses as small as one packet, the digital signal processor (DSP) interpolates what it thinks the audio should be and no problem is audible to the user.
- When jitter exceeds what the DSP can do to make up for the missing packets, audio problems are heard.



Note: In a properly designed network, packet loss should be near zero.



9.2 Traffic Characteristics

Traffic Characteristics

Network Traffic Trends

In the early 2000s, the predominant types of IP traffic were voice and data.

- Voice traffic has a predictable bandwidth need and known packet arrival times.
- Data traffic is not real-time and has unpredictable bandwidth need.
- Data traffic can temporarily burst, as when a large file is being downloaded. This bursting can consume the entire bandwidth of a link.

More recently, video traffic has become the increasingly important to business communications and operations.

- According to the Cisco Visual Networking Index (VNI), video traffic represented 70% of all traffic in 2017.
- By 2022, video will represent 82% of all traffic.
- Mobile video traffic will reach 60.9 exabytes per month by 2022.

The type of demands that voice, video, and data traffic place on the network are very different.

Voice traffic is predictable and smooth and very sensitive to delays and dropped packets.

- Voice packets must receive a higher priority than other types of traffic.
- Cisco products use the RTP port range 16384 to 32767 to prioritize voice traffic.

Voice can tolerate a certain amount of latency, jitter, and loss without any noticeable effects

Latency should be no more than 150 milliseconds (ms).

- Jitter should be no more than 30 ms, and packet loss no more than 1%.
- Voice traffic requires at least 30 Kbps of bandwidth.

Voice Traffic Characteristics	One-Way Requirements
<ul style="list-style-type: none">• Smooth• Benign• Drop sensitive• Delay sensitive• UDP priority	<ul style="list-style-type: none">• Latency \leq 150ms• Jitter \leq 30ms• Loss \leq 1% Bandwidth (30-128 Kbps)



Traffic Characteristics

Video

Video traffic tends to be unpredictable, inconsistent, and bursty. Compared to voice, video is less resilient to loss and has a higher volume of data per packet.

- The number and size of video packets varies every 33 ms based on the content of the video.
- UDP ports such as 554, are used for the Real-Time Streaming Protocol (RSTP) and should be given priority over other, less delay-sensitive, network traffic.
- Latency should be no more than 400 milliseconds (ms). Jitter should be no more than 50 ms, and video packet loss should be no more than 1%. Video traffic requires at least 384 Kbps of bandwidth.

Video Traffic Characteristics	One-Way Requirements
<ul style="list-style-type: none">• Bursty• Greedy• Drop sensitive• Delay sensitive• UDP priority	<ul style="list-style-type: none">• Latency \leq 200-400 ms• Jitter \leq 30-50 ms• Loss \leq 0.1 – 1%• Bandwidth (384 Kbps - 20 Mbps)



Traffic Characteristics Data

Data applications that have no tolerance for data loss, such as email and web pages, use TCP to ensure that if packets are lost in transit, they will be resent.

- Data traffic can be smooth or bursty.
- Network control traffic is usually smooth and predictable.

Some TCP applications can consume a large portion of network capacity. FTP will consume as much bandwidth as it can get when you download a large file, such as a movie or game.

Data Traffic Characteristics

- Smooth/bursty
- Benign/greedy
- Drop insensitive
- Delay insensitive
- TCP Retransmits



Traffic Characteristics Data (Cont.)

Data traffic is relatively insensitive to drops and delays compared to voice and video. Quality of Experience or QoE is important to consider with data traffic.

- Does the data come from an interactive application?
- Is the data mission critical?

Factor	Mission Critical	Not Mission Critical
Interactive	Prioritize for the lowest delay of all data traffic and strive for a 1 to 2 second response time.	Applications could benefit from lower delay.
Not interactive	Delay can vary greatly as long as the necessary minimum bandwidth is supplied.	Gets any leftover bandwidth after all voice, video, and other data application needs are met.



9.3 Queuing Algorithms



Queuing Algorithms

Queuing Overview

The QoS policy implemented by the network administrator becomes active when congestion occurs on the link. Queuing is a congestion management tool that can buffer, prioritize, and, if required, reorder packets before being transmitted to the destination.

A number of queuing algorithms are available:

- First-In, First-Out (FIFO)
- Weighted Fair Queuing (WFQ)
- Class-Based Weighted Fair Queuing (CBWFQ)
- Low Latency Queuing (LLQ)



Queuing Algorithms

First in First Out

- First In First Out (FIFO) queuing buffers and forwards packets in the order of their arrival.
- FIFO has no concept of priority or classes of traffic and consequently, makes no decision about packet priority.
- There is only one queue, and all packets are treated equally.
- Packets are sent out an interface in the order in which they arrive.

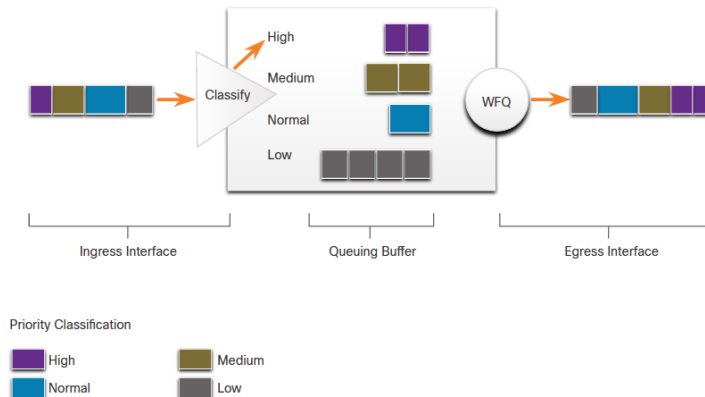


Queueing Algorithms

Weighted Fair Queuing (WFQ)

Weighted Fair Queuing (WFQ) is an automated scheduling method that provides fair bandwidth allocation to all network traffic.

- WFQ applies priority, or weights, to identified traffic, classifies it into conversations or flows, and then determines how much bandwidth each flow is allowed relative to other flows.
- WFQ classifies traffic into different flows based on source and destination IP addresses, MAC addresses, port numbers, protocol, and Type of Service (ToS) value.
- WFQ is not supported with tunneling and encryption because these features modify the packet content information required by WFQ for classification.





Class-Based Weighted Fair Queuing (CBWFQ)

Class-Based Weighted Fair Queuing (CBWFQ) extends the standard WFQ functionality to provide support for user-defined traffic classes.

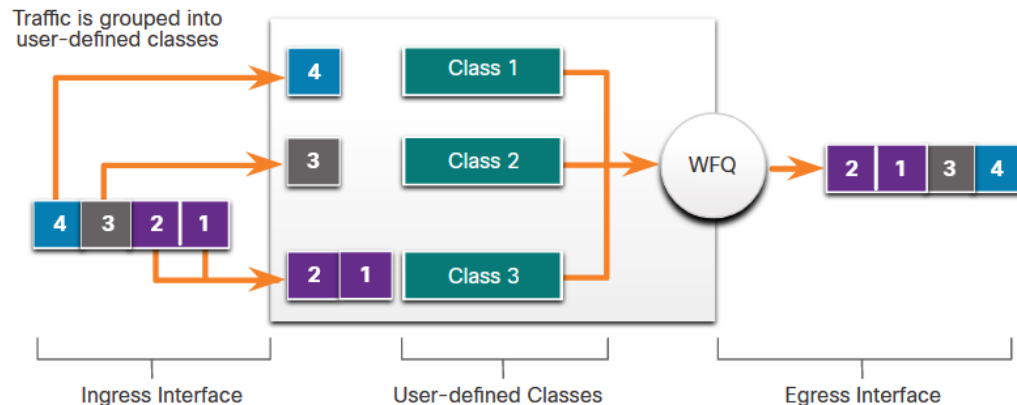
- Traffic classes are defined based on match criteria including protocols, access control lists (ACLs), and input interfaces.
- Packets satisfying the match criteria for a class constitute the traffic for that class.
- A FIFO queue is reserved for each class, and traffic belonging to a class is directed to the queue for that class.
- A class can be assigned characteristics, such as bandwidth, weight, and maximum packet limit. The bandwidth assigned to a class is the guaranteed bandwidth delivered during congestion.
- Packets belonging to a class are subject to the bandwidth and queue limits, which is the maximum number of packets allowed to accumulate in the queue, that characterize the class.



Class-Based Weighted Fair Queuing (CBWFQ) (Cont.)

After a queue has reached its configured queue limit, adding more packets to the class causes tail drop or packet drop to take effect, depending on how class policy is configured.

- Tail drop discards any packet that arrives at the tail end of a queue that has completely used up its packet-holding resources.
- This is the default queuing response to congestion. Tail drop treats all traffic equally and does not differentiate between classes of service.

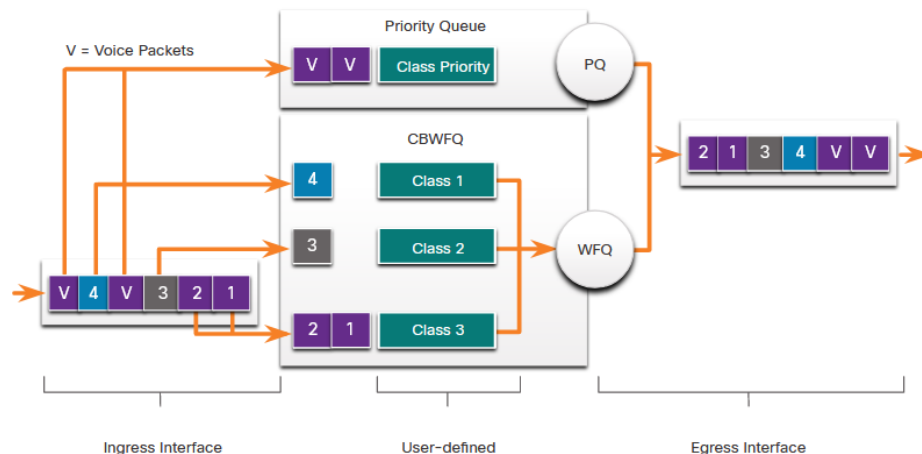


Queueing Algorithms

Low Latency Queuing (LLQ)

The Low Latency Queuing (LLQ) feature brings strict priority queuing (PQ) to CBWFQ.

- Strict PQ allows delay-sensitive packets such as voice to be sent before packets in other queues.
- LLQ allows delay-sensitive packets such as voice to be sent first (before packets in other queues), giving delay-sensitive packets preferential treatment over other traffic.
- Cisco recommends that only voice traffic be directed to the priority queue.





9.4 QoS Models



Selecting an Appropriate QoS Policy Model

There are three models for implementing QoS. QoS is implemented in a network using either IntServ or DiffServ.

- IntServ provides the highest guarantee of QoS, it is very resource-intensive, and therefore, not easily scalable.
- DiffServ is less resource-intensive and more scalable.
- IntServ and DiffServ are sometimes co-deployed in network QoS implementations.

Model	Description
Best-effort model	<ul style="list-style-type: none">• Not an implementation as QoS is not explicitly configured.• Use when QoS is not required.
Integrated services (IntServ)	<ul style="list-style-type: none">• Provides very high QoS to IP packets with guaranteed delivery.• Defines a signaling process for applications to signal to the network that they require special QoS for a period and that bandwidth should be reserved.• IntServ can severely limit the scalability of a network.
Differentiated services (DiffServ)	<ul style="list-style-type: none">• Provides high scalability and flexibility in implementing QoS.• Network devices recognize traffic classes and provide different levels of QoS to different traffic classes.

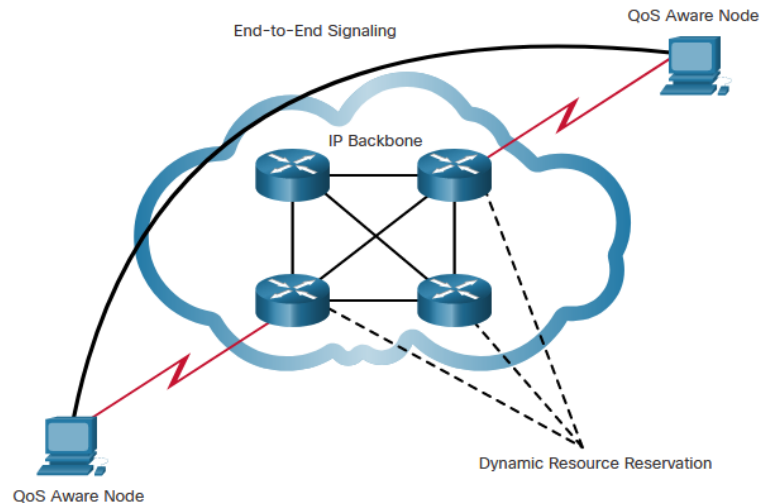
The basic design of the internet is best-effort packet delivery and provides no guarantees.

- The best-effort model treats all network packets in the same way, so an emergency voice message is treated the same way that a digital photograph attached to an email is treated.
- Benefits and drawbacks of the best effort model:

Benefits	Drawbacks
The model is the most scalable.	There are no guarantees of delivery.
Scalability is only limited by available bandwidth, in which case all traffic is equally affected.	Packets will arrive whenever they can and in any order possible, if they arrive at all.
No special QoS mechanisms are required.	No packets have preferential treatment.
It is the easiest and quickest model to deploy.	Critical data is treated the same as casual email is treated.

IntServ delivers the end-to-end QoS that real-time applications require.

- Explicitly manages network resources to provide QoS to individual flows or streams, sometimes called microflows.
- Uses resource reservation and admission-control mechanisms as building blocks to establish and maintain QoS.
- Uses a connection-oriented approach. Each individual communication must explicitly specify its traffic descriptor and requested resources to the network.
- The edge router performs admission control to ensure that available resources are sufficient in the network.



Integrated Services (Cont.)

In the IntServ model, the application requests a specific kind of service from the network before sending data.

- The application informs the network of its traffic profile and requests a particular kind of service that can encompass its bandwidth and delay requirements.
- IntServ uses the Resource Reservation Protocol (RSVP) to signal the QoS needs of an application's traffic along devices in the end-to-end path through the network.
- If network devices along the path can reserve the necessary bandwidth, the originating application can begin transmitting. If the requested reservation fails along the path, the originating application does not send any data.

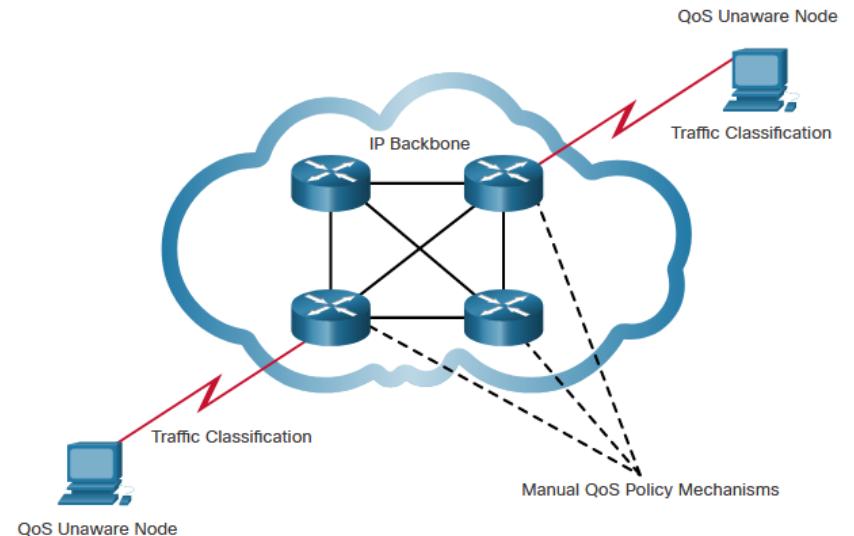
Benefits	Drawbacks
<ul style="list-style-type: none">• Explicit end-to-end resource admission control• Per-request policy admission control• Signaling of dynamic port numbers	<ul style="list-style-type: none">• Resource intensive due to the stateful architecture requirement for continuous signaling.• Flow-based approach not scalable to large implementations such as the internet.

QoS Models

Differentiated Services

The differentiated services (DiffServ) QoS model specifies a simple and scalable mechanism for classifying and managing network traffic.

- Is not an end-to-end QoS strategy because it cannot enforce end-to-end guarantees.
- Hosts forward traffic to a router which classifies the flows into aggregates (classes) and provides the appropriate QoS policy for the classes.
- Enforces and applies QoS mechanisms on a hop-by-hop basis, uniformly applying global meaning to each traffic class to provide both flexibility and scalability.





Differentiated Services (Cont.)

- DiffServ divides network traffic into classes based on business requirements. Each of the classes can then be assigned a different level of service.
- As the packets traverse a network, each of the network devices identifies the packet class and services the packet according to that class.
- It is possible to choose many levels of service with DiffServ.

Benefits	Drawbacks
<ul style="list-style-type: none">• Highly scalable• Provides many different levels of quality	<ul style="list-style-type: none">• No absolute guarantee of service quality• Requires a set of complex mechanisms to work in concert throughout the network



9.5 QoS Implementation Techniques

Packet loss is usually the result of congestion on an interface. Most applications that use TCP experience slowdown because TCP automatically adjusts to network congestion. Dropped TCP segments cause TCP sessions to reduce their window sizes. Some applications do not use TCP and cannot handle drops (fragile flows).

The following approaches can prevent drops in sensitive applications:

- Increase link capacity to ease or prevent congestion.
- Guarantee enough bandwidth and increase buffer space to accommodate bursts of traffic from fragile flows. WFQ, CBWFQ, and LLQ can guarantee bandwidth and provide prioritized forwarding to drop-sensitive applications.
- Drop lower-priority packets before congestion occurs. Cisco IOS QoS provides queuing mechanisms, such as weighted random early detection (WRED), that start dropping lower-priority packets before congestion occurs.

There are three categories of QoS tool, as described in the table.

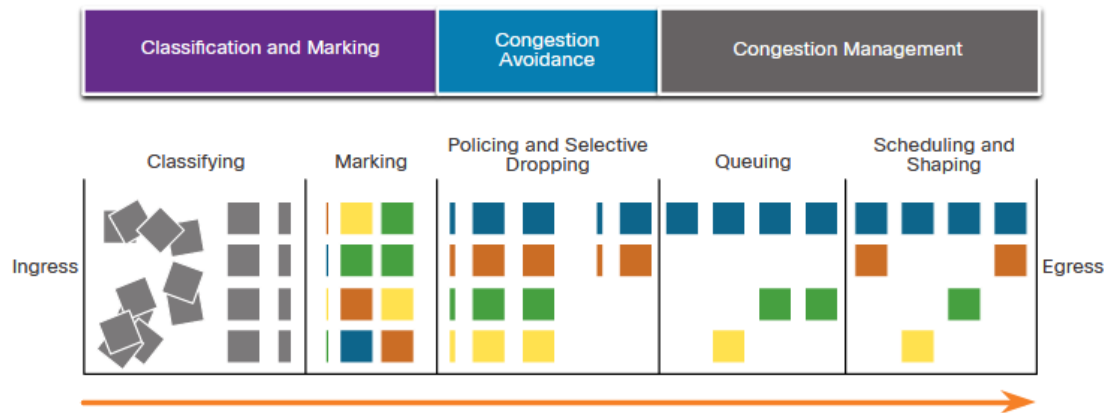
QoS Tools	Description
Classification and marking tools	<ul style="list-style-type: none">• Sessions, or flows, are analyzed to determine what traffic class they belong to.• When the traffic class is determined, the packets are marked.
Congestion avoidance tools	<ul style="list-style-type: none">• Traffic classes are allotted portions of network resources, as defined by the QoS policy.• The QoS policy also identifies how some traffic may be selectively dropped, delayed, or re-marked to avoid congestion.• The primary congestion avoidance tool is WRED and is used to regulate TCP data traffic in a bandwidth-efficient manner before tail drops caused by queue overflows occur.
Congestion management tools	<ul style="list-style-type: none">• When traffic exceeds available network resources, traffic is queued to await availability of resources.• Common Cisco IOS-based congestion management tools include CBWFQ and LLQ algorithms.

QoS Implementation Techniques

QoS Tools (Cont.)

The figure shows the sequence of QoS tools used when applied to packet flows.

- Ingress packets are classified and their respective IP header is marked.
- To avoid congestion, packets are then allocated resources based on defined policies.
- Packets are then queued and forwarded out the egress interface based on their defined QoS shaping and policing policy.



Note: Classification and marking can be done on ingress or egress, whereas other QoS actions such queuing and shaping are usually done on egress.

Classification and Marking

Before a packet can have a QoS policy applied to it, the packet has to be classified. Classification determines the class of traffic to which packets or frames belong. Only after traffic is marked can policies be applied to it.

How a packet is classified depends on the QoS implementation.

- Methods of classifying traffic flows at Layer 2 and 3 include using interfaces, ACLs, and class maps.
- Traffic can also be classified at Layers 4 to 7 using Network Based Application Recognition (NBAR).

Classification and Marking (Cont.)

How traffic is marked usually depends on the technology. The decision of whether to mark traffic at Layers 2 or 3 (or both) is not trivial and should be made after consideration of the following points:

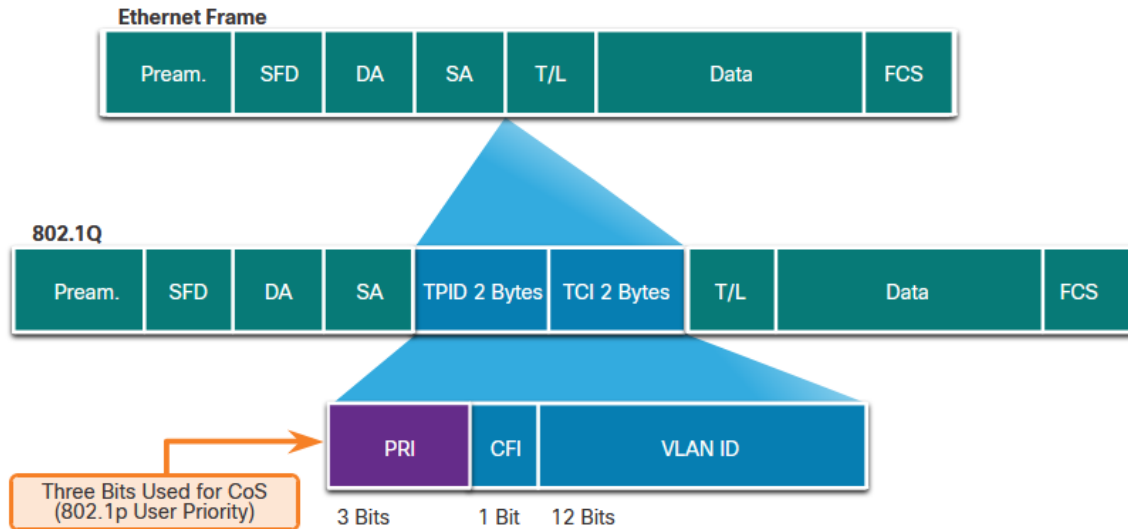
- Layer 2 marking of frames can be performed for non-IP traffic.
- Layer 2 marking of frames is the only QoS option available for switches that are not “IP aware”.
- Layer 3 marking will carry the QoS information end-to-end.

QoS Tools	Layer	Marking Field	Width in Bits
Ethernet (802.1q, 802.1p)	2	Class of Service (CoS)	3
802.11 (Wi-Fi)	2	Wi-Fi Traffic Identifier (TID)	3
MPLS	2	Experimental (EXP)	3
IPv4 and IPv6	3	IP Precedence (IPP)	3
IPv4 and IPv6	3	Differentiated Services Code Point (DSCP)	6

QoS Implementation Techniques

Marking at Layer 2

802.1Q is the IEEE standard that supports VLAN tagging at Layer 2 on Ethernet networks. When 802.1Q is implemented, two fields are inserted into the Ethernet frame following the source MAC address field.



Marking at Layer 2 (Cont.)

The 802.1Q standard also includes the QoS prioritization scheme known as IEEE 802.1p. The 802.1p standard uses the first three bits in the Tag Control Information (TCI) field. Known as the Priority (PRI) field, this 3-bit field identifies the Class of Service (CoS) markings.

Three bits means that a Layer 2 Ethernet frame can be marked with one of eight levels of priority (values 0-7).

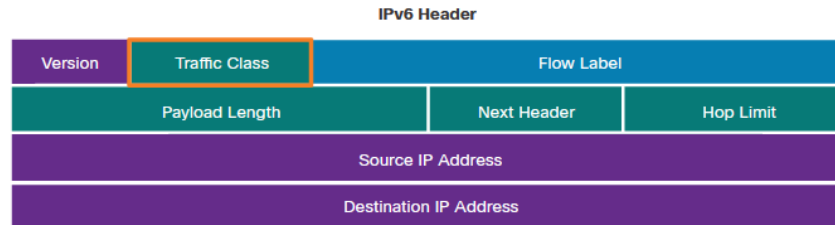
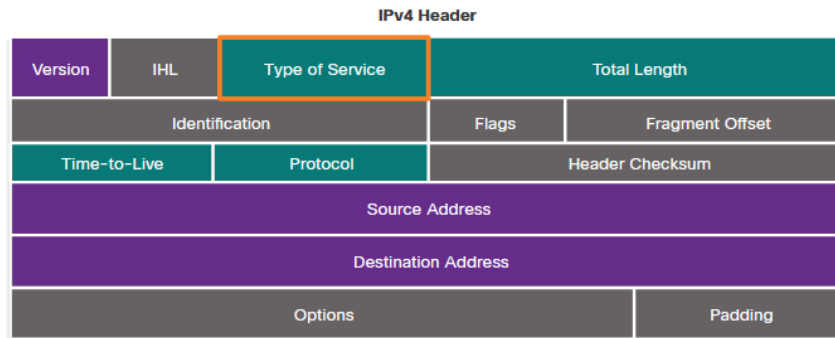
CoS Value	CoS Binary Value	Description
0	000	Best-Effort Data
1	001	Medium-Priority Data
2	010	High-Priority Data
3	011	Call Signaling
4	100	Videoconferencing
5	101	Voice bearer (voice traffic)
6	110	Reserved
7	111	Reserved

QoS Implementation Techniques

Marking at Layer 3

IPv4 and IPv6 specify an 8-bit field in their packet headers to mark packets.

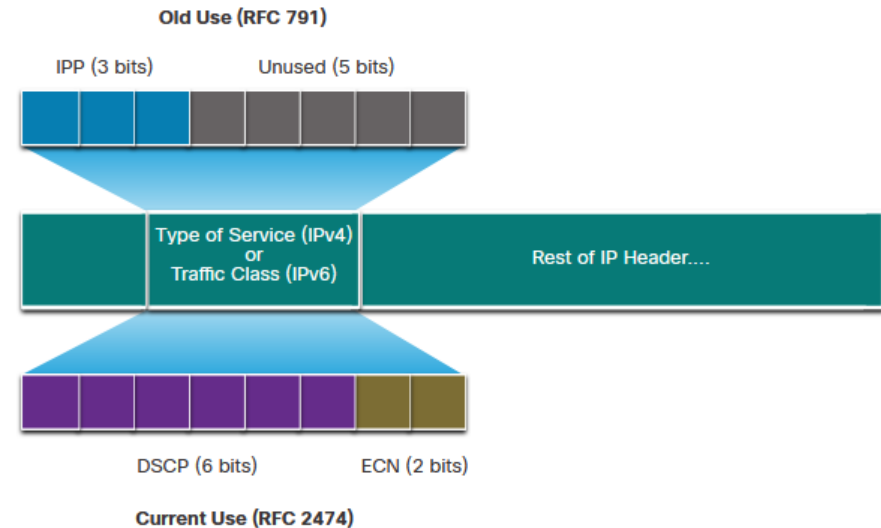
Both IPv4 and IPv6 support an 8-bit field for marking: the Type of Service (ToS) field for IPv4 and the Traffic Class field for IPv6.



Type of Service and Traffic Class Field

The Type of Service (IPv4) and Traffic Class (IPv6) carry the packet marking as assigned by the QoS classification tools.

- RFC 791 specified the 3-bit IP Precedence (IPP) field to be used for QoS markings.
- RFC 2474 supersedes RFC 791 and redefines the ToS field by renaming and extending the IPP field to 6 bits.
- Called the Differentiated Services Code Point (DSCP) field, these six bits offer a maximum of 64 possible classes of service.
- The remaining two IP Extended Congestion Notification (ECN) bits can be used by ECN-aware routers to mark packets instead of dropping them.



The 64 DSCP values are organized into three categories:

- **Best-Effort (BE)** - This is the default for all IP packets. The DSCP value is 0. The per-hop behavior is normal routing. When a router experiences congestion, these packets will be dropped. No QoS plan is implemented.
- **Expedited Forwarding (EF)** - RFC 3246 defines EF as the DSCP decimal value 46 (binary **101110**). The first 3 bits (101) map directly to the Layer 2 CoS value 5 used for voice traffic. At Layer 3, Cisco recommends that EF only be used to mark voice packets.
- **Assured Forwarding (AF)** - RFC 2597 defines AF to use the 5 most significant DSCP bits to indicate queues and drop preference.

Assured Forwarding values are shown in the figure.

The **AFxy** formula is specified as follows:

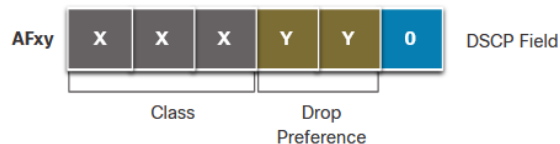
- The first 3 most significant bits are used to designate the class. Class 4 is the best queue and Class 1 is the worst queue.
- The 4th and 5th most significant bits are used to designate the drop preference.
- The 6th most significant bit is set to zero.

Best Queue



Worst Queue

Assured Forwarding Values			
	Low Drop	Medium Drop	High Drop
Class 4	AF41 (34)	AF42 (36)	AF43 (38)
Class 3	AF31 (26)	AF32 (28)	AF33 (30)
Class 2	AF21 (18)	AF22 (20)	AF23 (22)
Class 1	AF11 (10)	AF12 (12)	AF13 (14)



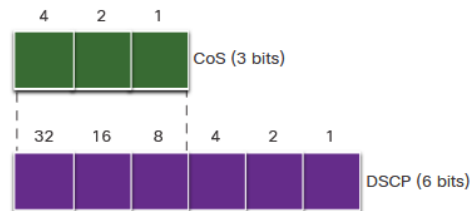
For example: AF32 belongs to class 3 (binary 011) and has a medium drop preference (binary 10). The full DSCP value is 28 because you include the 6th 0 bit (binary 011100).

QoS Implementation Techniques

Class Selector Bits

Class Selector (CS) bits:

- The first 3 most significant bits of the DSCP field and indicate the class.
- Map directly to the 3 bits of the CoS field and the IPP field to maintain compatibility with 802.1p and RFC 791.



CoS values, Class Selectors, and corresponding DSCP 6-bit value

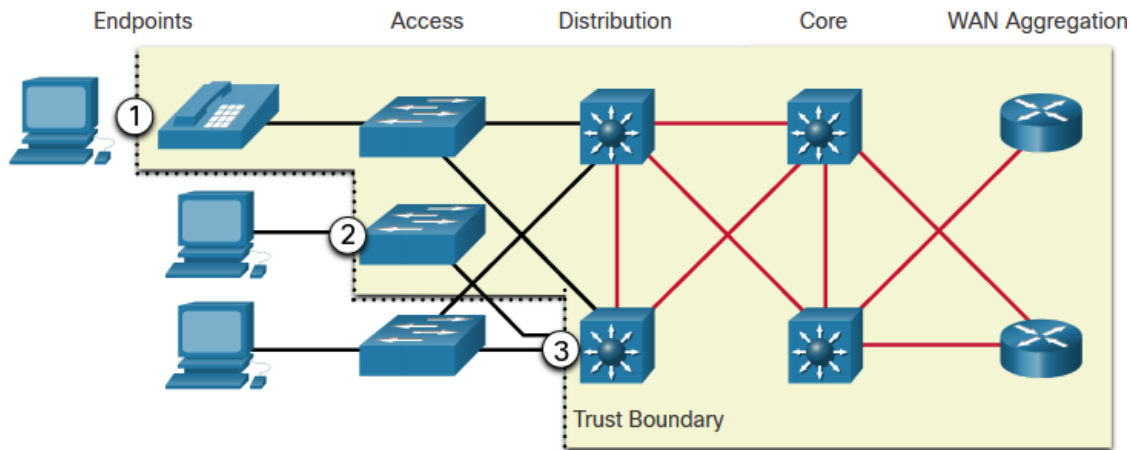
CoS Value	CoS Binary Value	Class Selector (CS)	CS Binary	DSCP Decimal Value
0	000	CS0*/DF	000 000	0
1	001	CS1	001 000	8
2	010	CS2	010 000	16
3	011	CS3	011 000	24
4	100	CS4	100 000	32
5	101	CS5	101 000	40
6	110	CS6	110 000	48
7	111	CS7	111 000	56

QoS Implementation Techniques

Trust Boundaries

Traffic should be classified and marked as close to its source as technically and administratively feasible. This defines the trust boundary.

1. Trusted endpoints have the capabilities and intelligence to mark application traffic to the appropriate Layer 2 CoS and/or Layer 3 DSCP values.
2. Secure endpoints can have traffic marked at the Layer 2 switch.
3. Traffic can also be marked at Layer 3 switches / routers.



Congestion avoidance tools monitor network traffic loads in an effort to anticipate and avoid congestion at common network and internetwork bottlenecks before congestion becomes a problem.

- They monitor network traffic loads in an effort to anticipate and avoid congestion at common network and internetwork bottlenecks before congestion becomes a problem.
- They monitor the average depth of the queue. When the queue is below the minimum threshold, there are no drops. As the queue fills up to the maximum threshold, a small percentage of packets are dropped. When the maximum threshold is passed, all packets are dropped.

Some congestion avoidance techniques provide preferential treatment for which packets get dropped.

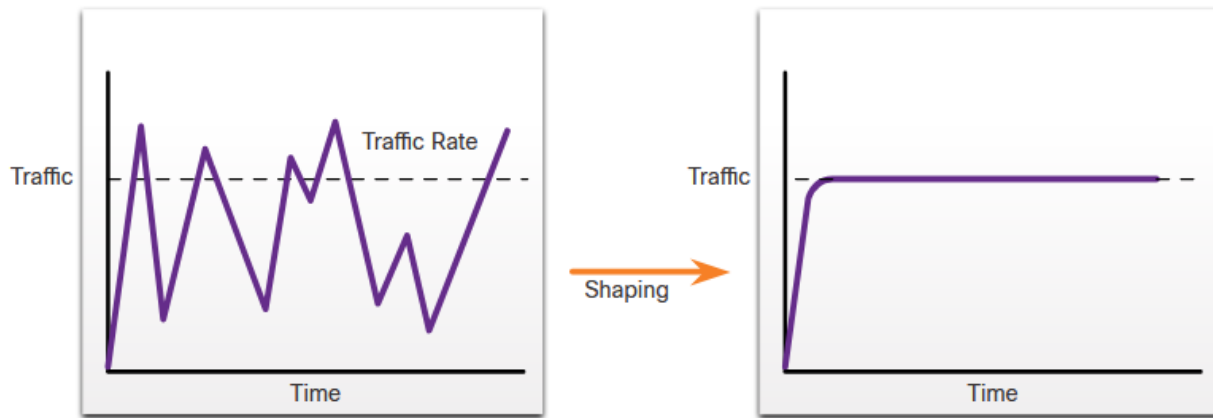
- Weighted random early detection (WRED) allows for congestion avoidance on network interfaces by providing buffer management and allowing TCP traffic to decrease, or throttle back, before buffers are exhausted.
- WRED helps avoid tail drops and maximizes network use and TCP-based application performance.

QoS Implementation Techniques

Shaping and Policing

Traffic shaping and traffic policing are two mechanisms provided by Cisco IOS QoS software to prevent congestion.

- Traffic shaping retains excess packets in a queue and then schedules the excess for later transmission over increments of time. Traffic shaping results in a smoothed packet output rate.
- Shaping is an outbound concept; packets going out an interface get queued and can be shaped. In contrast, policing is applied to inbound traffic on an interface.



QoS Implementation Techniques

Shaping and Policing (Cont.)

Policing is applied to inbound traffic on an interface. Policing is commonly implemented by service providers to enforce a contracted customer information rate (CIR). However, the service provider may also allow bursting over the CIR if the service provider's network is not currently experiencing congestion.



QoS policies must consider the full path from source to destination.

A few guidelines that help ensure the best experience for end users includes the following:

- Enable queuing at every device in the path between source and destination.
- Classify and mark traffic as close the source as possible.
- Shape and police traffic flows as close to their sources as possible.

