



Traffic Shaping

Cisco IOS Quality of Service (QoS) software includes four types of traffic shaping:

- Generic Traffic Shaping (GTS)
- Class-Based Traffic Shaping
- Frame Relay Traffic Shaping (FRTS)
- Distributed Traffic Shaping (DTS)

All four traffic shaping methods are similar in implementation, though their command line interfaces (CLIs) differ somewhat and they use different types of queues to contain and shape traffic that is deferred. If a packet is deferred, GTS and Class-Based Shaping use a weighted fair queue to hold the delayed traffic. FRTS uses either a custom queue or a priority queue.

This section explains how traffic shaping works and describes the Cisco IOS QoS traffic shaping mechanisms. It also described traffic-shaping feature called Low Latency Traffic Shaping (LLQ).

This section contains the following topics:

- [About Traffic Shaping, page 10-2](#)
- [Generic Traffic Shaping, page 10-6](#)
- [Class-Based Traffic Shaping, page 10-8](#)
- [Frame Relay Traffic Shaping, page 10-9](#)
- [Distributed Traffic Shaping, page 10-12](#)
- [Low-Latency Queueing, page 10-14](#)

About Traffic Shaping

Traffic shaping allows you to control outgoing traffic on an interface to match the traffic speed of the remote target interface and to ensure that the traffic conforms to specific policies. Traffic that adheres to a particular profile can be shaped to meet downstream requirements, thereby eliminating bottlenecks in topologies caused by data-rate mismatches.

Why Use Traffic Shaping?

The primary reasons to use traffic shaping are to control access to available bandwidth, to ensure that traffic conforms to specific policies, and to regulate the flow of traffic in order to avoid congestion. Some example reasons for using traffic shaping follow:

- Control access to bandwidth when policy dictates that the average rate of a given interface should not exceed a certain rate.
- Configure traffic shaping on an interface if you have a network with differing access rates. Suppose that one end of the link in a Frame Relay network runs at 256 kbps and the other end of the link runs at 128 kbps. Sending packets at 256 kbps could cause failure of the applications using the link.

A similar, more complicated case would be a link-layer network giving indications of congestion with differing access rates on different attached data terminal equipment (DTE) devices. The network may be able to deliver more transit speed to a given DTE device at a specific time than at another time.

- If you offer a subrate service, traffic shaping enables you to use the router to partition your T1 or T3 links into smaller channels.

Traffic shaping prevents packet loss. Its use is especially important in Frame Relay networks because the switch cannot determine which packets take precedence or which packets should be dropped when congestion occurs.

Traffic Shaping and Rate of Transfer

Traffic shaping limits the rate of transmission of data. You can limit the data transfer to one of the following:

- A specific configured rate
- A derived rate based on the level of congestion

The rate of transfer depends on three components that constitute the token bucket: burst size, mean rate, and measurement (time) interval. The mean rate is equal to the burst size divided by the interval.

When traffic shaping is enabled, the bit rate of the interface does not exceed the mean rate over any integral multiple of the interval. During every interval, the burst size is usually the maximum number of bits that can be sent. Within the interval, however, the bit rate may be faster than the mean rate at any given time.

One additional variable applies to traffic shaping: Excess Burst Size (called the *Be size*). The *Be Size* corresponds to the number of noncommitted bits—those outside the committed information rate (CIR)—that are still accepted by the Frame Relay switch but are marked as discard eligible (DE).

The *Be size* allows more than the burst size to be sent during a time interval in certain situations. The switch allows the packets belonging to the Excess Burst to go through but it will mark them by setting the DE bit. Whether the packets are sent depends on how the switch is configured.

When the *Be size* equals 0, the interface sends no more than the burst size every interval, achieving an average rate no higher than the mean rate. However, when the *Be size* is greater than 0, the interface can send as many as $B_c + B_e$ bits in a burst, if the maximum amount was not sent in a previous time period. Whenever the number of bits sent during an interval is less than the burst size, the remaining number of bits can be sent in a later interval.

Discard Eligible Bit

You can specify which Frame Relay packets have low priority or low time sensitivity. These packets are the first to be dropped when a Frame Relay switch is congested. The Discard Eligible (DE) bit allows a Frame Relay switch to identify such packets.

You can define DE lists that identify the characteristics of packets to be eligible for discarding, and you can also specify DE groups to identify the data-link connection identifier (DLCI) that is affected.

You can specify DE lists based on the protocol or the interface. You can also specify DE lists that are based on characteristics such as fragmentation of the packet, a specific TCP or User Datagram Protocol (UDP) port, an access list number, or a packet size.

Differences Between Shaping Mechanisms

GTS, Class-Based Shaping, DTS, and FRTS are similar in implementation, sharing the same code and data structures, but they differ in regard to their CLIs and the queue types they use.

Here are a few ways in which these mechanisms differ:

- For GTS, the shaping queue is a weighted fair queue. For FRTS, the queue can be a weighted fair queue (configured by the **frame-relay fair-queue** command), a strict priority queue with weighted fair queueing (WFQ) (configured by the **frame-relay ip rtp priority** command in addition to the **frame-relay fair-queue** command), custom queueing (CQ), priority queueing (PQ), or first-in, first-out queueing (FIFO).
- For Class-Based Shaping, GTS can be configured on a class, rather than only on an access control list (ACL). You must first define traffic classes based on match criteria including protocols, ACLs, and input interfaces. You can then apply traffic shaping to each defined class.
- FRTS supports shaping on a per-DLCI basis; GTS and DTS are configurable per interface or subinterface.
- DTS supports traffic shaping based on a variety of match criteria, including user-defined classes, and differentiated services code point (DSCP).

Table 10-1 summarizes these differences.

Table 10-1 Differences Between Shaping Mechanisms

Mechanism	GTS	Class-Based	DTS	FRTS
Command-Line Interface	<ul style="list-style-type: none"> Applies parameters per subinterface traffic group command supported 	<ul style="list-style-type: none"> Applies parameters per interface or per class 	<ul style="list-style-type: none"> Applies parameters per interface or subinterface 	<ul style="list-style-type: none"> Classes of parameters Applies parameters to all virtual circuits (VCs) on an interface through inheritance mechanism No traffic group command
Queues Supported	<ul style="list-style-type: none"> WFQ per subinterface 	<ul style="list-style-type: none"> class-based weighted fair queuing (CBWFQ) inside GTS 	<ul style="list-style-type: none"> WFQ, strict priority queue with WFQ, CQ, PQ, first-come, first-served (FCFS) per VC 	<ul style="list-style-type: none"> WFQ, strict priority queue with WFQ, CQ, PQ, FCFS per VC

You can configure GTS to behave the same as FRTS by allocating one DLCI per subinterface and using GTS plus backward explicit congestion notification (BECN) support.

Traffic Shaping and Queueing

Traffic shaping smooths traffic by storing traffic above the configured rate in a queue.

When a packet arrives at the interface for transmission, the following sequence occurs:

1. If the queue is empty, the arriving packet is processed by the traffic shaper.
 - If possible, the traffic shaper sends the packet.
 - Otherwise, the packet is placed in the queue.
2. If the queue is not empty, the packet is placed in the queue.

When packets are in the queue, the traffic shaper removes the number of packets it can send from the queue at each time interval.

Generic Traffic Shaping

Generic Traffic Shaping (GTS) shapes traffic by reducing outbound traffic flow to avoid congestion. GTS constrains traffic to a particular bit rate using the token bucket mechanism. See the section “What is a Token Bucket” in the *Cisco IOS Quality of Service Solutions Configuration Guide, Release 12.2*:

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgr/fqos_c/fqcp4/qcfcpolsh.htm#1000909

How It Works

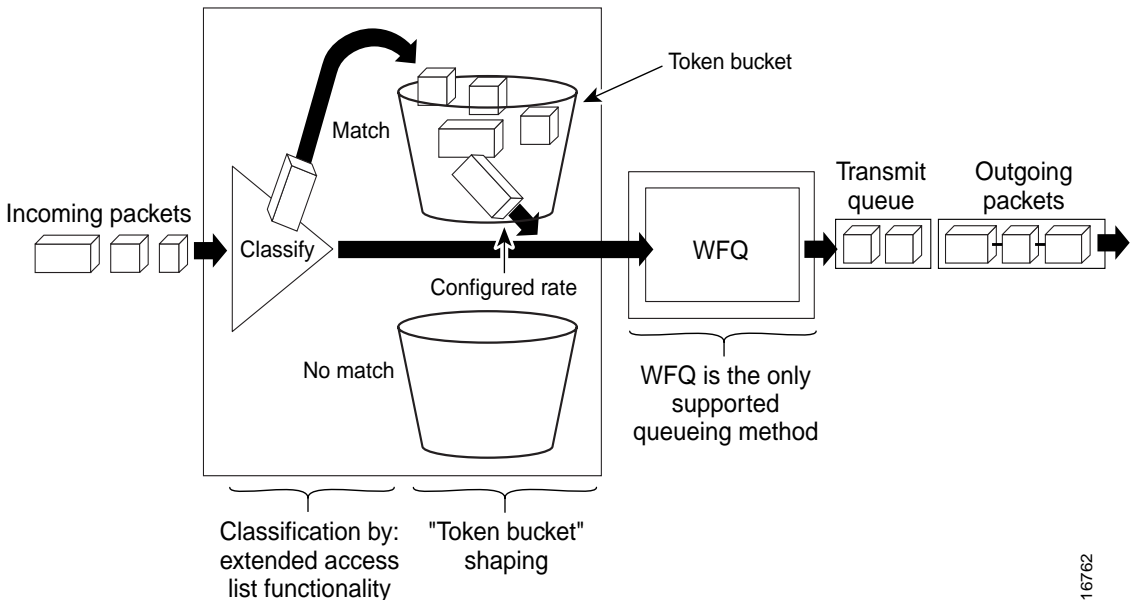
GTS applies traffic shaping on a per-interface basis and can use access lists to select the traffic to shape. GTS works with a variety of Layer 2 technologies, including Frame Relay, ATM, Switched Multimegabit Data Service (SMDS), and Ethernet.

On a Frame Relay subinterface, GTS can be set up to adapt dynamically to available bandwidth by integrating backward explicit congestion notification (BECN) signals. GTS also can be shape traffic to a specified rate. GTS can be configured on an ATM/ATM Interface Processor (AIP) interface to respond to the Resource Reservation Protocol (RSVP) feature signalled over statically configured ATM permanent virtual circuits (PVCs).

GTS is supported on most media and encapsulation types on the router. GTS can be applied to a specific access list on an interface.

Figure 10-1 shows how GTS works.

Figure 10-1 Generic Traffic Shaping



16762

Configuration and Commands

For information on how to configure GTS, see the chapter “Configuring Generic Traffic 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/fqcp4/qcfigts.htm#80560

For information on traffic shaping commands, see the *Cisco IOS Quality of Service Solutions Command Reference, Release 12.2* manual:

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

Class-Based Traffic Shaping

Class-Based Traffic Shaping can be enabled on any interface that supports GTS.

How It Works

Using Class-Based Traffic Shaping, you can perform the following tasks:

- Configure GTS on a traffic class to provide greater flexibility for configuring traffic shaping. Previously, this ability was limited to the use of ACLs.
- Specify average rate or peak rate traffic shaping. This type of shaping allows more data than the CIR to be sent if the bandwidth is available.
- Configure class-based weighted fair queueing (CBWFQ) inside GTS. CBWFQ allows you to specify the exact amount of bandwidth to allocate for a specific class of traffic. You can configure up to 64 classes and control their distribution.

Flow-based WFQ applies weights to traffic to classify the traffic into conversations and determine how much bandwidth each conversation is allowed. These weights and traffic classifications are dependent on and limited to the seven IP Precedence levels.

CBWFQ allows you to define what constitutes a class based on criteria that exceed the confines of flow. CBWFQ allows you to use ACLs and protocols or input interface names to define how traffic will be classified, thereby providing coarser granularity. You need not maintain traffic classification on a flow basis. Moreover, you can configure up to 64 discrete classes in a service policy.

Configuration and Commands

For information on how to configure Class-Based Shaping, see the chapter “Configuring Class-Based 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/qcfcshp.htm#80464

For information on traffic shaping commands, see the *Cisco IOS Quality of Service Solutions Command Reference, Release 12.2* manual:

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

Restrictions

Peak and average traffic shaping is configured on a per-interface or per-class basis, and cannot be used in conjunction with commands used to configure GTS from previous versions of Cisco IOS. These commands include the following:

- **traffic-shape adaptive**
- **traffic-shape fecn-adaptive**
- **traffic-shape group**
- **traffic-shape rate**

Adaptive traffic shaping for Frame Relay networks is not supported using the Class-Based Shaping feature. To configure adaptive GTS for Frame Relay networks, you must use the commands from releases prior to Release 12.1(2) of Cisco IOS software.

Frame Relay Traffic Shaping

Cisco has long provided support for FECN for DECnet and OSI, and BECN for Systems Network Architecture (SNA) traffic using Logical Link Control, type 2 (LLC2) encapsulation via RFC 1490 and DE bit support. FRTS builds upon this existing Frame Relay support with additional capabilities that improve the scalability and performance of a Frame Relay network, increasing the density of VCs and improving response time.

As is also true of GTS, FRTS can eliminate bottlenecks in Frame Relay networks that have high-speed connections at the central site and low-speed connections at branch sites. You can configure rate enforcement—a peak rate configured to limit outbound traffic—to limit the rate at which data is sent on the VC at the central site.

How It Works

Using FRTS, you can configure rate enforcement to either the CIR or some other defined value such as the excess information rate on a per-VC basis. The ability to allow the transmission speed used by the router to be controlled by criteria other than line speed provides a mechanism for sharing media by multiple VCs. You can allocate bandwidth to each VC, creating a virtual time-division multiplexing (TDM) network.

You can also define PQ, CQ, and WFQ at the VC or subinterface level. Using these queueing methods allows for finer granularity in the prioritization and queueing of traffic, providing more control over the traffic flow on an individual VC. If you combine CQ with the per-VC queueing and rate enforcement capabilities, you enable Frame Relay VCs to carry multiple traffic types such as IP, SNA, and Internetwork Packet Exchange (IPX) with bandwidth guaranteed for each traffic type.

Using information contained in the BECN-tagged packets received from the network, FRTS can also dynamically throttle traffic. With BECN-based throttling, packets are held in the buffers of the router to reduce the data flow from the router into the Frame Relay network. The throttling is done on a per-VC basis and the transmission rate is adjusted based on the number of BECN-tagged packets received.

With the Cisco FRTS feature, you can integrate ATM ForeSight closed-loop congestion control to actively adapt to downstream congestion conditions.

Derived Rates

In Frame Relay networks, BECNs and FECNs indicate congestion. BECN and FECN are specified by bits within a Frame Relay frame.

FECNs are generated when data is sent out a congested interface; they indicate to a DTE device that congestion was encountered. Traffic is marked with BECN if the queue for the opposite direction is deep enough to trigger FECNs at the current time.

BECNs notify the sender to decrease the transmission rate. If the traffic is one-way only (such as multicast traffic), there is no reverse traffic with BECNs to notify the sender to slow down. Thus, when a DTE device receives an FECN, it first determines if it is sending any data in return. If it is sending return data, this

data will get marked with a BECN on its way to the other DTE device. However, if the DTE device is not sending any data, the DTE device can send a Q.922 TEST RESPONSE message with the BECN bit set.

When an interface configured with traffic shaping receives a BECN, it immediately decreases its maximum rate by a large amount. If, after several intervals, the interface has not received another BECN and traffic is waiting in the queue, the maximum rate increases slightly. The dynamically adjusted maximum rate is called the derived rate.

The derived rate will always be between the upper bound and the lower bound configured on the interface.

Configuration and Commands

For more information on configuring Frame Relay, refer to the chapter “Configuring Frame Relay” in the *Cisco IOS Wide-Area Networking Configuration Guide, Release 12.2* manual:

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

For information on configuring Frame Relay as it relates to voice traffic, refer to the chapter “Configuring Voice Over Frame Relay” in the *Cisco IOS Voice, Video, and Fax Configuration Guide, Release 12.2* manual:

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

For information on Frame Relay commands, see the section “Frame Relay Commands” in the *Cisco IOS Wide-Area Networking Command Reference, Release 12.2* manual:

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

For information on Frame Relay ATM commands, see the section “Frame Relay—ATM Internetworking Commands” in the *Cisco IOS Wide-Area Networking Command Reference, Release 12.2* manual:

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

Restrictions

FRTS applies only to Frame Relay PVCs and switched virtual circuits (SVCs).

Distributed Traffic Shaping

DTS provides a method of managing the bandwidth of an interface to avoid congestion, to meet remote site requirements, and to conform to a service rate that is provided on that interface.

DTS uses queues to buffer traffic surges that can congest a network and send the data to the network at a regulated rate. This ensures that traffic will behave to the configured descriptor, as defined by the CIR, Bc, and Be. With the defined average bit rate and burst size that is acceptable on that shaped entity, you can derive a time interval value.

Prerequisites

Distributed Cisco Express Forwarding (dCEF) must be enabled on the interface before DTS can be enabled.

A policy map and class maps must be created before DTS is enabled.

How It Works

The Be size allows more than the Bc size to be sent during a time interval under certain conditions. Therefore, DTS provides two types of **shape** commands: **average** and **peak**. When **shape average** is configured, the interface sends no more than the Bc size for each interval, achieving an average rate no higher than the CIR. When the **shape peak** command is configured, the interface sends Bc plus Be bits in each interval.

In a link layer network such as Frame Relay, the network sends messages with the forward explicit congestion notification (FECN) or BECN if there is congestion. With the DTS feature, the traffic shaping adaptive mode takes advantage of these signals and adjusts the traffic descriptors, thereby regulating the amount of traffic entering or leaving the interface accordingly.

DTS provides the following key benefits:

- Offloads traffic shaping from the Route Switch Processor (RSP) to the VIP.
- Supports up to 200 shape queues per VIP, supporting up to OC-3 rates when the average packet size is 250 bytes or greater and when using a VIP2-50 or better with eight MB of SRAM. Line rates below T3 are supported with a VIP2-40.
- Configures DTS at the interface level or subinterface level.
- Shaping based on the following traffic match criteria:
 - Access list
 - Packet marking
 - Input port
 - Other matching criteria. For information about other matching criteria, see the section “Creating a Traffic Class” in the *Cisco IOS Quality of Service Solutions Command Reference, Release 12.2* manual.
- Optional configuration to respond to Frame Relay network congestion by reducing the shaped-to rate for a period of time until congestion is believed to have subsided. Supports FECN, BECN, and ForeSight Frame Relay signaling.

This feature runs on Cisco 7500 series routers with VIP2-40, VIP2-50, or greater.

Configuration

For information on how to configure DTS, see the chapter “Configuring Distributed Traffic Shaping” in the *Cisco IOS Quality of Service Solutions Command Reference, Release 12.2* manual:

http://www.cisco.com/univercd/cc/td/doc/product/software/ios122/122cgcr/fqos_c/fqcp4/qcfdts.htm

Restrictions

DTS does not support the following:

- Fast EtherChannel, Multilink PPP (MLP), Tunnel, VLANs, and dialer interface
- Any VIP below a VIP2-40



Note

A VIP2-50 is strongly recommended when the aggregate line rate of the port adapters on the VIP is greater than DS3. A VIP2-50 card is required for OC-3 rates.

Low-Latency Queueing

The LLQ feature brings strict PQ to CBWFQ. Strict PQ allows delay-sensitive data such as voice to be dequeued and sent before packets in other queues are dequeued.

Without LLQ, CBWFQ provides WFQ based on defined classes with no strict priority queue available for real-time traffic. CBWFQ allows you to define traffic classes and then assign characteristics to these classes. For example, you can designate the minimum bandwidth delivered to the class during congestion.

For CBWFQ, the weight for a packet belonging to a specific class is derived from the bandwidth you assigned to the class at configuration. Therefore, the bandwidth assigned to the packets of a class determines the order in which packets are sent. All packets are serviced fairly based on weight; no class of packets may be granted strict priority. This scheme poses problems for voice traffic that is largely intolerant of delay, especially variation in delay. For voice traffic, variations in delay introduce irregularities of transmission, such as *jitter* in the heard conversation.

LLQ provides strict priority queueing for CBWFQ, reducing jitter in voice conversations. Configured by the priority command, LLQ enables use of a single, strict priority queue within CBWFQ at the class level. This allows you to direct traffic belonging to a class to the CBWFQ strict priority queue. To enqueue class traffic to the strict priority queue, specify the named class within a policy map and then configure the priority command for the class. (Classes to which the priority command is applied are considered priority classes.) Within a policy map, you can

give one or more classes priority status. When multiple classes within a single policy map are configured as priority classes, all traffic from these classes is enqueued to the same single, strict priority queue.

One of the ways in which the strict PQ used within CBWFQ differs from its use outside CBWFQ is in the parameters it takes. Outside CBWFQ, you can use the **ip rtp priority** command to specify the range of UDP ports whose voice traffic flows are given priority service.

Using the priority command, you are no longer limited to a UDP port number to stipulate priority flows because you can configure the priority status for a class within CBWFQ. Instead, all of the valid match criteria used to specify traffic for a class now apply to priority traffic. These methods of specifying traffic for a class include matching on access lists, protocols, and input interfaces. Moreover, within an access list you can specify that traffic matches are allowed based on the IP differentiated services code point (DSCP) value. This value is set using the first six bits of the ToS byte in the IP header.

Although it is possible to enqueue various types of real-time traffic to the strict priority queue, Cisco strongly recommends that you direct only voice traffic to this queue. The reason is that voice traffic is well-behaved, whereas other types of real-time traffic are not well-behaved. Moreover, voice traffic requires nonvariable delays to avoid jitter.

Real-time traffic such as video could introduce variation in delay, thereby thwarting the steadiness of delay required for successful voice traffic transmission.

For more conceptual information about LLQ, see the section “Weighted Fair Queueing” 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#xtocid46014

For information on how to configure LLQ, see the chapter “Configuring Weighted Fair Queueing” 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/qcfwfq.htm

