



CHAPTER 57

Configuring QoS

Have you ever participated in a long-distance phone call that involved a satellite connection? The conversation might be interrupted with brief, but perceptible, gaps at odd intervals. Those gaps are the time, called the latency, between the arrival of packets being transmitted over the network. Some network traffic, such as voice and video, cannot tolerate long latency times. Quality of service (QoS) is a feature that lets you give priority to critical traffic, prevent bandwidth hogging, and manage network bottlenecks to prevent packet drops.

This chapter describes how to apply QoS policies and includes the following sections:

- [Information About QoS, page 57-1](#)
- [Licensing Requirements for QoS, page 57-5](#)
- [Guidelines and Limitations, page 57-5](#)
- [Configuring QoS, page 57-6](#)
- [Feature History for QoS, page 57-11](#)

Information About QoS

You should consider that in an ever-changing network environment, QoS is not a one-time deployment, but an ongoing, essential part of network design.



Note

QoS is only available in single context mode.

This section describes the QoS features supported by the security appliance and includes the following topics:

- [Supported QoS Features, page 57-2](#)
- [What is a Token Bucket?, page 57-2](#)
- [Information About Policing, page 57-3](#)
- [Information About Priority Queuing, page 57-3](#)
- [Information About Traffic Shaping, page 57-4](#)
- [DSCP and DiffServ Preservation, page 57-5](#)

Supported QoS Features

The security appliance supports the following QoS features:

- Policing—To prevent individual flows from hogging the network bandwidth, you can limit the maximum bandwidth used per flow. See the “[Information About Policing](#)” section on page 57-3 for more information.
- Priority queuing—For critical traffic that cannot tolerate latency, such as Voice over IP (VoIP), you can identify traffic for Low Latency Queuing (LLQ) so that it is always transmitted ahead of other traffic. See the “[Information About Priority Queuing](#)” section on page 57-3 for more information.
- Traffic shaping—If you have a device that transmits packets at a high speed, such as a security appliance with Fast Ethernet, and it is connected to a low speed device such as a cable modem, then the cable modem is a bottleneck at which packets are frequently dropped. To manage networks with differing line speeds, you can configure the security appliance to transmit packets at a fixed slower rate. See the “[Information About Traffic Shaping](#)” section on page 57-4 for more information.

What is a Token Bucket?

A token bucket is used to manage a device that regulates the data in a flow. For example, the regulator might be a traffic policer or a traffic shaper. A token bucket itself has no discard or priority policy. Rather, a token bucket discards tokens and leaves to the flow the problem of managing its transmission queue if the flow overdrives the regulator.

A token bucket is a formal definition of a rate of transfer. It has three components: a burst size, an average rate, and a time interval. Although the average rate is generally represented as bits per second, any two values may be derived from the third by the relation shown as follows:

average rate = burst size / time interval

Here are some definitions of these terms:

- Average rate—Also called the committed information rate (CIR), it specifies how much data can be sent or forwarded per unit time on average.
- Burst size—Also called the Committed Burst (Bc) size, it specifies in bits or bytes per burst how much traffic can be sent within a given unit of time to not create scheduling concerns. (For traffic shaping, it specifies bits per burst; for policing, it specifies bytes per burst.)
- Time interval—Also called the measurement interval, it specifies the time quantum in seconds per burst.

In the token bucket metaphor, tokens are put into the bucket at a certain rate. The bucket itself has a specified capacity. If the bucket fills to capacity, newly arriving tokens are discarded. Each token is permission for the source to send a certain number of bits into the network. To send a packet, the regulator must remove from the bucket a number of tokens equal in representation to the packet size.

If not enough tokens are in the bucket to send a packet, the packet either waits until the bucket has enough tokens (in the case of traffic shaping) or the packet is discarded or marked down (in the case of policing). If the bucket is already full of tokens, incoming tokens overflow and are not available to future packets. Thus, at any time, the largest burst a source can send into the network is roughly proportional to the size of the bucket.

Note that the token bucket mechanism used for traffic shaping has both a token bucket and a data buffer, or queue; if it did not have a data buffer, it would be a policer. For traffic shaping, packets that arrive that cannot be sent immediately are delayed in the data buffer.

For traffic shaping, a token bucket permits burstiness but bounds it. It guarantees that the burstiness is bounded so that the flow will never send faster than the token bucket capacity, divided by the time interval, plus the established rate at which tokens are placed in the token bucket. See the following formula:

$$(\text{token bucket capacity in bits} / \text{time interval in seconds}) + \text{established rate in bps} = \text{maximum flow speed in bps}$$

This method of bounding burstiness also guarantees that the long-term transmission rate will not exceed the established rate at which tokens are placed in the bucket.

Information About Policing

Policing is a way of ensuring that no traffic exceeds the maximum rate (in bits/second) that you configure, thus ensuring that no one traffic flow or class can take over the entire resource. When traffic exceeds the maximum rate, the security appliance drops the excess traffic. Policing also sets the largest single burst of traffic allowed.

Information About Priority Queuing

LLQ priority queuing lets you prioritize certain traffic flows (such as latency-sensitive traffic like voice and video) ahead of other traffic.

The security appliance supports two types of priority queuing:

- Standard priority queuing—Standard priority queuing uses an LLQ priority queue on an interface (see the [“Configuring the Standard Priority Queue for an Interface”](#) section on page 57-7), while all other traffic goes into the “best effort” queue. Because queues are not of infinite size, they can fill and overflow. When a queue is full, any additional packets cannot get into the queue and are dropped. This is called *tail drop*. To avoid having the queue fill up, you can increase the queue buffer size. You can also fine-tune the maximum number of packets allowed into the transmit queue. These options let you control the latency and robustness of the priority queuing. Packets in the LLQ queue are always transmitted before packets in the best effort queue.
- Hierarchical priority queuing—Hierarchical priority queuing is used on interfaces on which you enable a traffic shaping queue. A subset of the shaped traffic can be prioritized. The standard priority queue is not used. See the following guidelines about hierarchical priority queuing:
 - Priority packets are always queued at the head of the shape queue so they are always transmitted ahead of other non-priority queued packets.
 - Priority packets are never dropped from the shape queue unless the sustained rate of priority traffic exceeds the shape rate.
 - For IPsec-encrypted packets, you can only match traffic based on the DSCP or precedence setting.
 - IPsec-over-TCP is not supported for priority traffic classification.

Information About Traffic Shaping

Traffic shaping is used to match device and link speeds, thereby controlling packet loss, variable delay, and link saturation, which can cause jitter and delay.

**Note**

Traffic shaping is not supported on the ASA 5580.

- Traffic shaping must be applied to all outgoing traffic on a physical interface or in the case of the ASA 5505, on a VLAN. You cannot configure traffic shaping for specific types of traffic.
- Traffic shaping is implemented when packets are ready to be transmitted on an interface, so the rate calculation is performed based on the actual size of a packet to be transmitted, including all the possible overhead such as the IPsec header and L2 header.
- The shaped traffic includes both through-the-box and from-the-box traffic.
- The shape rate calculation is based on the standard token bucket algorithm. The token bucket size is twice the Burst Size value. See the [“What is a Token Bucket?”](#) section on page 57-2.
- When bursty traffic exceeds the specified shape rate, packets are queued and transmitted later. Following are some characteristics regarding the shape queue (for information about hierarchical priority queuing, see the [“Information About Priority Queuing”](#) section on page 57-3):
 - The queue size is calculated based on the shape rate. The queue can hold the equivalent of 200-milliseconds worth of shape rate traffic, assuming a 1500-byte packet. The minimum queue size is 64.
 - When the queue limit is reached, packets are tail-dropped.
 - Certain critical keep-alive packets such as OSPF Hello packets are never dropped.
 - The time interval is derived by $time_interval = burst_size / average_rate$. The larger the time interval is, the burstier the shaped traffic might be, and the longer the link might be idle. The effect can be best understood using the following exaggerated example:

Average Rate = 1000000

Burst Size = 1000000

In the above example, the time interval is 1 second, which means, 1 Mbps of traffic can be bursted out within the first 10 milliseconds of the 1-second interval on a 100 Mbps FE link and leave the remaining 990 milliseconds idle without being able to send any packets until the next time interval. So if there is delay-sensitive traffic such as voice traffic, the Burst Size should be reduced compared to the average rate so the time interval is reduced.

How QoS Features Interact

You can configure each of the QoS features alone if desired for the security appliance. Often, though, you configure multiple QoS features on the security appliance so you can prioritize some traffic, for example, and prevent other traffic from causing bandwidth problems.

See the following supported feature combinations per interface:

- Standard priority queuing (for specific traffic) + Policing (for the rest of the traffic).
You cannot configure priority queuing and policing for the same set of traffic.
- Traffic shaping (for all traffic on an interface) + Hierarchical priority queuing (for a subset of traffic).

You cannot configure traffic shaping and standard priority queuing for the same interface; only hierarchical priority queuing is allowed. For example, if you configure standard priority queuing for the global policy, and then configure traffic shaping for a specific interface, the feature you configured last is rejected because the global policy overlaps the interface policy.

Typically, if you enable traffic shaping, you do not also enable policing for the same traffic, although the security appliance does not restrict you from configuring this.

DSCP and DiffServ Preservation

- DSCP markings are preserved on all traffic passing through the security appliance.
- The security appliance does not locally mark/remark any classified traffic, but it honors the Expedited Forwarding (EF) DSCP bits of every packet to determine if it requires “priority” handling and will direct those packets to the LLQ.
- DiffServ marking is preserved on packets when they traverse the service provider backbone so that QoS can be applied in transit (QoS tunnel pre-classification).

Licensing Requirements for QoS

The following table shows the licensing requirements for this feature:

Model	License Requirement
All models	Base License.

Guidelines and Limitations

This section includes the guidelines and limitations for this feature.

Context Mode Guidelines

Supported in single context mode only. Does not support multiple context mode.

Firewall Mode Guidelines

Supported in routed firewall mode only. Does not support transparent firewall mode.

IPv6 Guidelines

Does not support IPv6.

Model Guidelines

Traffic shaping is not supported on the ASA 5580.

Additional Guidelines and Limitations

- For traffic shaping, you can only use the **class-default** class map, which is automatically created by the security appliance, and which matches all traffic.

- For priority traffic, you cannot use the **class-default** class map.
- For hierarchical priority queuing, for encrypted VPN traffic, you can only match traffic based on the DSCP or precedence setting; you cannot match a tunnel group.
- For hierarchical priority queuing, IPsec-over-TCP traffic is not supported.
- You cannot configure traffic shaping and standard priority queuing for the same interface; only hierarchical priority queuing is allowed.
- For standard priority queuing, the queue must be configured for a physical interface or for a VLAN on the ASA 5505.
- You cannot create a standard priority queue for a Ten Gigabit Ethernet interface; priority queuing is not necessary for an interface with high bandwidth.

Configuring QoS

This section includes the following topics:

- [Determining the Queue and TX Ring Limits for a Standard Priority Queue, page 57-6](#)
- [Configuring the Standard Priority Queue for an Interface, page 57-7](#)
- [Configuring a Service Rule for Standard Priority Queuing and Policing, page 57-8](#)
- [Configuring a Service Rule for Traffic Shaping and Hierarchical Priority Queuing, page 57-9](#)

Determining the Queue and TX Ring Limits for a Standard Priority Queue

To determine the priority queue and TX ring limits, use the worksheets below.

Table 57-1 shows how to calculate the priority queue size. Because queues are not of infinite size, they can fill and overflow. When a queue is full, any additional packets cannot get into the queue and are dropped (called *tail drop*). To avoid having the queue fill up, you can adjust the queue buffer size according to the “[Configuring the Standard Priority Queue for an Interface](#)” section on page 57-7.

Table 57-1 Queue Limit Worksheet

Step 1	$\frac{\text{_____ Mbps}}{\text{_____}} \times 125 = \text{_____}$ <p style="text-align: center;"><i>Outbound bandwidth</i> <i>(Mbps or Kbps)¹</i></p> $\frac{\text{_____ Kbps}}{\text{_____}} \times .125 = \text{_____}$ <p style="text-align: center;"><i># of bytes/ms</i></p>
Step 2	$\frac{\text{_____}}{\text{_____}} \div \frac{\text{_____}}{\text{_____}} \times \text{_____} = \text{_____}$ <p style="text-align: center;"><i># of bytes/ms from Step 1</i> <i>Average packet size (bytes)²</i> <i>Delay (ms)³</i> Queue limit (# of packets)</p>

1. For example, DSL might have an uplink speed of 768 Kbps. Check with your provider.
2. Determine this value from a codec or sampling size. For example, for VoIP over VPN, you might use 160 bytes. We recommend 256 bytes if you do not know what size to use.
3. The delay depends on your application. For example, the recommended maximum delay for VoIP is 200 ms. We recommend 500 ms if you do not know what delay to use.

Table 57-2 shows how to calculate the TX ring limit. This limit determines the maximum number of packets allowed into the Ethernet transmit driver before the driver pushes back to the queues on the interface to let them buffer packets until the congestion clears. This setting guarantees that the hardware-based transmit ring imposes a limited amount of extra latency for a high-priority packet.

Table 57-2 TX Ring Limit Worksheet

Step 1	$\frac{\text{Outbound bandwidth (Mbps or Kbps)}^1}{\text{Mbps}} \times 125 = \text{\# of bytes/ms}$ $\frac{\text{Outbound bandwidth (Mbps or Kbps)}^1}{\text{Kbps}} \times 0.125 = \text{\# of bytes/ms}$				
Step 2	$\frac{\text{\# of bytes/ms from Step 1}}{\text{Maximum packet size (bytes)}^2} \times \text{Delay (ms)}^3 =$				$\text{TX ring limit (\# of packets)}$

1. For example, DSL might have an uplink speed of 768 Kbps. Check with your provider.
2. Typically, the maximum size is 1538 bytes, or 1542 bytes for tagged Ethernet. If you allow jumbo frames (if supported for your platform), then the packet size might be larger.
3. The delay depends on your application. For example, to control jitter for VoIP, you should use 20 ms.

Configuring the Standard Priority Queue for an Interface

If you enable standard priority queuing for traffic on a physical interface, then you need to also create the priority queue on each interface. Each physical interface uses two queues: one for priority traffic, and the other for all other traffic. For the other traffic, you can optionally configure policing.



Note

The standard priority queue is not required for hierarchical priority queuing with traffic shaping; see the “[Information About Priority Queuing](#)” section on page 57-3 for more information.

Restrictions

You cannot create a priority queue for a Ten Gigabit Ethernet interface; priority queuing is not necessary for an interface with high bandwidth.

Detailed Steps

- Step 1** Go to Configuration > Device Management > Advanced > Priority Queue, and click **Add**. The Add Priority Queue dialog box displays.
- Step 2** From the Interface drop-down list, choose the physical interface name on which you want to enable the priority queue, or for the ASA 5505, the VLAN interface name.

- Step 3** To change the size of the priority queues, in the Queue Limit field, enter the number of average, 256-byte packets that the specified interface can transmit in a 500-ms interval.
- A packet that stays more than 500 ms in a network node might trigger a timeout in the end-to-end application. Such a packet can be discarded in each network node.
- Because queues are not of infinite size, they can fill and overflow. When a queue is full, any additional packets cannot get into the queue and are dropped (called *tail drop*). To avoid having the queue fill up, you can use this option to increase the queue buffer size.
- The upper limit of the range of values for this option is determined dynamically at run time. The key determinants are the memory needed to support the queues and the memory available on the device.
- The Queue Limit that you specify affects both the higher priority low-latency queue and the best effort queue.
- Step 4** To specify the depth of the priority queues, in the Transmission Ring Limit field, enter the number of maximum 1550-byte packets that the specified interface can transmit in a 10-ms interval.
- This setting guarantees that the hardware-based transmit ring imposes no more than 10-ms of extra latency for a high-priority packet.
- This option sets the maximum number of low-latency or normal priority packets allowed into the Ethernet transmit driver before the driver pushes back to the queues on the interface to let them buffer packets until the congestion clears.
- The upper limit of the range of values is determined dynamically at run time. The key determinants are the memory needed to support the queues and the memory available on the device.
- The Transmission Ring Limit that you specify affects both the higher priority low-latency queue and the best-effort queue.
-

Configuring a Service Rule for Standard Priority Queuing and Policing

You can configure standard priority queuing and policing for different class maps within the same policy map. See the [“How QoS Features Interact” section on page 57-4](#) for information about valid QoS configurations.

To create a policy map, perform the following steps.

Restrictions

- You cannot use the **class-default** class map for priority traffic.
- You cannot configure traffic shaping and standard priority queuing for the same interface; only hierarchical priority queuing is allowed.

Guidelines

- For priority traffic, identify only latency-sensitive traffic.
- For policing traffic, you can choose to police all other traffic, or you can limit the traffic to certain types.

Detailed Steps

-
- Step 1** To configure priority queuing, configure a service policy rule in the Configuration > Firewall > Service Policy Rules pane according to [Chapter 23, “Configuring Service Policy Rules.”](#)
- You can configure QoS as part of a new service policy rule, or you can edit an existing service policy.
- Step 2** In the Rule Actions dialog box, click the **QoS** tab.
- Step 3** Click **Enable priority for this flow**.
- If this service policy rule is for an individual interface, ASDM automatically creates the priority queue for the interface (Configuration > Device Management > Advanced > Priority Queue; for more information, see the [“Configuring the Standard Priority Queue for an Interface”](#) section on page 57-7). If this rule is for the global policy, then you need to manually add the priority queue to one or more interfaces *before* you configure the service policy rule.
- Step 4** Click **Finish**. The service policy rule is added to the rule table.
- Step 5** To configure policing, configure a service policy rule for the same interface in the Configuration > Firewall > Service Policy Rules pane according to [Chapter 23, “Configuring Service Policy Rules.”](#)
- For policing traffic, you can choose to police all traffic that you are not prioritizing, or you can limit the traffic to certain types.
- Step 6** In the Rule Actions dialog box, click the **QoS** tab.
- Step 7** Click **Enable policing**, then check the **Input policing** or **Output policing** (or both) check boxes to enable the specified type of traffic policing. For each type of traffic policing, configure the following fields:
- **Committed Rate**—The rate limit for this traffic flow; this is a value in the range 8000-2000000000, specifying the maximum speed (bits per second) allowed.
 - **Conform Action**—The action to take when the rate is less than the conform-burst value. Values are transmit or drop.
 - **Exceed Action**—Take this action when the rate is between the conform-rate value and the conform-burst value. Values are transmit or drop.
 - **Burst Rate**—A value in the range 1000-512000000, specifying the maximum number of instantaneous bytes allowed in a sustained burst before throttling to the conforming rate value.
- Step 8** Click **Finish**. The service policy rule is added to the rule table.
- Step 9** Click **Apply** to send the configuration to the device.
-

Configuring a Service Rule for Traffic Shaping and Hierarchical Priority Queuing

You can configure traffic shaping for all traffic on an interface, and optionally hierarchical priority queuing for a subset of latency-sensitive traffic.

Guidelines

- One side-effect of priority queuing is packet re-ordering. For IPsec packets, out-of-order packets that are not within the anti-replay window generate warning syslog messages. These warnings are false alarms in the case of priority queuing. You can configure the IPsec anti-replay window size to avoid possible false alarms. See the Configuration > VPN > IPsec > IPsec Rules > Enable Anti-replay window size option in the [“Adding Crypto Maps”](#) section on page 67-12.

- For hierarchical priority queuing, you do not need to create a priority queue on an interface.

Restrictions

- For hierarchical priority queuing, for encrypted VPN traffic, you can only match traffic based on the DSCP or precedence setting; you cannot match a tunnel group.
- For hierarchical priority queuing, IPsec-over-TCP traffic is not supported.
- Traffic shaping is not supported on the ASA 5580.
- For traffic shaping, you can only use the **class-default** class map, which is automatically created by the security appliance, and which matches all traffic.
- You cannot configure traffic shaping and standard priority queuing for the same interface; only hierarchical priority queuing is allowed. See the [“How QoS Features Interact”](#) section on page 57-4 for information about valid QoS configurations.
- You cannot configure traffic shaping in the global policy.

Detailed Steps

-
- Step 1** Configure a service policy on the Configuration > Firewall > Service Policy Rules pane according to [Chapter 23, “Configuring Service Policy Rules.”](#)
- You can configure QoS as part of a new service policy rule, or you can edit an existing service policy.
- Step 2** In the Rule Actions dialog box, click the **QoS** tab.
- Step 3** Click **Enable traffic shaping**, and configure the following fields:
- **Average Rate**—Sets the average rate of traffic in bits per second over a given fixed time period, between 64000 and 154400000. Specify a value that is a multiple of 8000.
 - **Burst Size**—Sets the average burst size in bits that can be transmitted over a given fixed time period, between 2048 and 154400000. Specify a value that is a multiple of 128. If you do not specify the Burst Size, the default value is equivalent to 4-milliseconds of traffic at the specified Average Rate. For example, if the average rate is 1000000 bits per second, 4 ms worth = $1000000 * 4/1000 = 4000$.
- Step 4** (Optional) To configure priority queuing for a subset of shaped traffic:
- a. Click **Enforce priority to selected shape traffic**.
 - b. Click **Configure** to identify the traffic that you want to prioritize.
You are prompted to identify the traffic for which you want to apply priority queuing.
 - c. After you identify the traffic (see the [“Adding a Service Policy Rule for Through Traffic”](#) section on page 23-7), click **Next**.
 - d. Click **Enable priority for this flow**.
 - e. Click **Finish**.
You return to the QoS tab.
- Step 5** Click **Finish**. The service policy rule is added to the rule table.
- Step 6** Click **Apply** to send the configuration to the device.
-

Feature History for QoS

Table 57-3 lists each feature change and the platform release in which it was implemented. ASDM is backwards-compatible with multiple platform releases, so the specific ASDM release in which support was added is not listed.

Table 57-3 Feature History for QoS

Feature Name	Platform Releases	Feature Information
Priority queuing and policing	7.0(1)	We introduced QoS priority queuing and policing. We introduced the following screens: Configuration > Device Management > Advanced > Priority Queue Configuration > Firewall > Service Policy Rules
Shaping and hierarchical priority queuing	7.2(4)/8.0(4)	We introduced QoS shaping and hierarchical priority queuing. We modified the following screen: Configuration > Firewall > Service Policy Rules.

