

Traffic Management User Guide for QFabric Systems

Published
2023-12-14

Juniper Networks, Inc.
1133 Innovation Way
Sunnyvale, California 94089
USA
408-745-2000
www.juniper.net

Juniper Networks, the Juniper Networks logo, Juniper, and Junos are registered trademarks of Juniper Networks, Inc. in the United States and other countries. All other trademarks, service marks, registered marks, or registered service marks are the property of their respective owners.

Juniper Networks assumes no responsibility for any inaccuracies in this document. Juniper Networks reserves the right to change, modify, transfer, or otherwise revise this publication without notice.

Traffic Management User Guide for QFabric Systems
Copyright © 2023 Juniper Networks, Inc. All rights reserved.

The information in this document is current as of the date on the title page.

YEAR 2000 NOTICE

Juniper Networks hardware and software products are Year 2000 compliant. Junos OS has no known time-related limitations through the year 2038. However, the NTP application is known to have some difficulty in the year 2036.

END USER LICENSE AGREEMENT

The Juniper Networks product that is the subject of this technical documentation consists of (or is intended for use with) Juniper Networks software. Use of such software is subject to the terms and conditions of the End User License Agreement ("EULA") posted at <https://support.juniper.net/support/eula/>. By downloading, installing or using such software, you agree to the terms and conditions of that EULA.

Table of Contents

[About This Guide | xi](#)

1

[CoS Overview](#)

[Overview | 2](#)

[Overview of Junos OS CoS | 2](#)

[Overview of Policers | 5](#)

[Basic Concepts | 15](#)

[Configuring CoS | 15](#)

[Understanding Junos CoS Components | 22](#)

[Assigning CoS Components to Interfaces | 27](#)

[Understanding CoS Packet Flow | 29](#)

[Understanding Default CoS Settings | 33](#)

[CoS Inputs and Outputs Overview | 47](#)

[CoS Upgrade Requirements and Feature Change Compatibility | 49](#)

[Overview of CoS Upgrade Requirements \(Junos OS Release 11.1 or 11.2 to a Later Release\) | 49](#)

[Overview of CoS Upgrade Requirements to Junos OS Release 12.2 | 51](#)

[Overview of CoS Upgrade Requirements to Junos OS Release 12.3 \(QFX3500 and QFX3600 Switches\) or to Junos OS Release 13.1 \(QFabric Systems\) | 53](#)

[Overview of CoS Changes Introduced in Junos OS Release 11.3 | 57](#)

[Overview of CoS Changes Introduced in Junos OS Release 12.2 | 67](#)

2

[Classifying Traffic \(Classifiers, Forwarding Classes, and Rewrite Rules\)](#)

[Using Classifiers, Forwarding Classes, and Rewrite Rules | 71](#)

[Understanding CoS Classifiers | 72](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 81](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 84](#)

Example: Configuring Unicast Classifiers | 86

Requirements | 87

Overview | 88

Verification | 89

Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers | 91

Requirements | 92

Overview | 92

Verification | 93

Understanding Host Inbound Traffic Classification | 94

Understanding Default CoS Scheduling and Classification | 95

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 106

Understanding CoS Code-Point Aliases | 120

Defining CoS Code-Point Aliases | 123

Understanding CoS Forwarding Classes | 124

Defining CoS Forwarding Classes | 131

Example: Configuring Forwarding Classes | 133

Requirements | 134

Overview | 134

Example 1: Configuring Forwarding Classes for Switches Except QFX10000 | 136

Verification | 137

Example 2: Configuring Forwarding Classes for QFX10000 Switches | 138

Verification | 139

Understanding CoS Forwarding Class Sets (Priority Groups) | 140

Defining CoS Forwarding Class Sets | 142

Example: Configuring Forwarding Class Sets | 143

Requirements | 144

Overview | 144

Verification | 146

Understanding Host Routing Engine Outbound Traffic Queues and Defaults | 148

Changing the Host Outbound Traffic Default Queue Mapping | 151

Understanding CoS Rewrite Rules | 152

Defining CoS Rewrite Rules | 154

Configuring CoS Fixed Classifier Rewrite Values for Native FC Interfaces (NP_Ports) | 156

Troubleshooting an Unexpected Rewrite Value | 159

Scheduling Traffic

Using Schedulers (Node Devices) | 163

Understanding Default CoS Scheduling and Classification | 164

Understanding CoS Scheduling on QFabric System Node Device Fabric (fte) Ports | 174

Understanding CoS Scheduling Behavior and Configuration Considerations | 179

Understanding CoS Output Queue Schedulers | 186

Defining CoS Queue Schedulers | 194

Example: Configuring Queue Schedulers | 198

Requirements | 200

Overview | 200

Verification | 203

Defining CoS Queue Scheduling Priority | 206

Example: Configuring Queue Scheduling Priority | 207

Requirements | 209

Overview | 209

Verification | 211

Understanding CoS Traffic Control Profiles | 213

Understanding CoS Priority Group Scheduling | 214

Defining CoS Traffic Control Profiles (Priority Group Scheduling) | 218

Example: Configuring Traffic Control Profiles (Priority Group Scheduling) | 220

Requirements | 221

Overview | 221

Verification | 222

Understanding CoS Hierarchical Port Scheduling (ETS) | 223

Example: Configuring CoS Hierarchical Port Scheduling (ETS) | 230

Requirements | 231

Overview | 232

Configuration | 238

Verification | 252

Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth | 266

Understanding CoS Priority Group Shaping and Queue Shaping (Maximum Bandwidth) | 269

Example: Configuring Minimum Guaranteed Output Bandwidth | 272

Requirements | 273

Overview | 274

Verification | 276

Troubleshooting Egress Bandwidth That Exceeds the Configured Minimum Bandwidth | 279

Example: Configuring Maximum Output Bandwidth | 280

Requirements | 282

Overview | 283

Verification | 284

Troubleshooting Egress Bandwidth That Exceeds the Configured Maximum Bandwidth | 287

Troubleshooting Egress Queue Bandwidth Impacted by Congestion | 288

Understanding CoS WRED Drop Profiles | 290

Configuring CoS WRED Drop Profiles | 297

Drop Profiles on Switches Except QFX10000 | 298

Drop Profiles on QFX 10000 Switches | 299

Example: Configuring WRED Drop Profiles | 300

Requirements | 301

Overview | 301

Configuring WRED Drop Profiles on Switches Except QFX10000 | 302

Verification | 304

Configuring WRED Drop Profiles on QFX10000 Switches | 305

Verification | 306

Configuring CoS Drop Profile Maps | 307

Example: Configuring Drop Profile Maps | 307

Requirements | 309

Overview | 309

Verification | 309

Troubleshooting a Port Reset on QFabric Systems When a Queue Stops Transmitting Traffic | 311

Using Schedulers (Interconnect Device Fabric) | 315

Understanding Default CoS Scheduling on QFabric System Interconnect Devices (Junos OS Release 13.1 and Later Releases) | 315

Understanding CoS Scheduling Across the QFabric System | 327

Example: Configuring CoS Scheduling Across the QFabric System | 353

Requirements | 353

Overview | 353

Configuration | 365

Verification | 381

Understanding CoS Fabric Forwarding Class Sets | 396

Configuring CoS Fabric Forwarding Class Set Scheduler Maps (Fabric Scheduler to Fabric FC-Set Mapping) | 409

Understanding How to Mitigate Fate Sharing on a QFabric System Interconnect Device by Remapping Traffic Flows (Forwarding Classes) | 410

Configuring Fate Sharing Mitigation Across the Interconnect Device by Remapping Traffic Flows (Forwarding Classes) | 434

4

Configuring Data Center Bridging (ETS, PFC, DCBX) and Flow Control

Using Data Center Bridging and Flow Control | 443

Understanding DCB Features and Requirements | 444

Understanding CoS Hierarchical Port Scheduling (ETS) | 447

Example: Configuring CoS Hierarchical Port Scheduling (ETS) | 454

Requirements | 455

Overview | 456

Configuration | 462

Verification | 476

Disabling the ETS Recommendation TLV | 490

Understanding CoS Flow Control (Ethernet PAUSE and PFC) | **491**

Enabling and Disabling CoS Symmetric Ethernet PAUSE Flow Control | **504**

Configuring CoS Asymmetric Ethernet PAUSE Flow Control | **505**

Configuring CoS PFC (Congestion Notification Profiles) | **507**

Example: Configuring CoS PFC for FCoE Traffic | **510**

Requirements | **511**

Overview | **511**

Configuration | **514**

Verification | **521**

Troubleshooting Dropped FCoE Traffic | **524**

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | **528**

Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch) | **549**

Requirements | **550**

Overview | **550**

Configuration | **553**

Verification | **556**

Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface | **562**

Requirements | **562**

Overview | **563**

Configuration | **566**

Verification | **569**

Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces | **575**

Requirements | **575**

Overview | **576**

Configuration | **581**

Verification | **586**

Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI) | **594**

Requirements | **595**

- Overview | 595
- Configuration | 602
- Verification | 610

Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway | 620

Example: Configuring IEEE 802.1p Priority Remapping on an FCoE-FC Gateway | 624

- Requirements | 624
- Overview | 625
- Configuration | 629
- Verification | 632

Understanding DCBX | 638

Configuring the DCBX Mode | 648

Configuring DCBX Autonegotiation | 649

Understanding DCBX Application Protocol TLV Exchange | 652

Defining an Application for DCBX Application Protocol TLV Exchange | 657

Configuring an Application Map for DCBX Application Protocol TLV Exchange | 658

Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange | 660

Example: Configuring DCBX Application Protocol TLV Exchange | 661

- Requirements | 662
- Overview | 662
- Configuration | 667
- Verification | 670

5

Configuring Buffers

Using Buffers | 677

Understanding CoS Buffer Configuration | 677

Configuring Global Ingress and Egress Shared Buffers | 701

Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic | 703

- Requirements | 704
- Overview | 704
- Configuration | 706

Verification | 709

Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled | 712

Requirements | 713

Overview | 713

Configuration | 715

Verification | 718

Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 720

Requirements | 721

Overview | 721

Configuration | 723

Verification | 726

Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 729

Requirements | 730

Overview | 730

Configuration | 732

Verification | 735

6

Learn About Technology

Learn About Technology | 739

Data Center Technology Overview Videos | 739

7

Configuration Statements and Operational Commands

Monitoring Interfaces That Have CoS Components | 742

Monitoring CoS Classifiers | 744

Monitoring CoS Forwarding Classes | 746

Monitoring CoS Rewrite Rules | 750

Monitoring CoS Code-Point Value Aliases | 752

Monitoring CoS Scheduler Maps | 754

Junos CLI Reference Overview | 756

About This Guide

Use this guide to understand and configure class of service (CoS) features in Junos OS to define service levels that provide different delay, jitter, and packet loss characteristics to particular applications served by specific traffic flows. Applying CoS features to each device in your network ensures quality of service (QoS) for traffic throughout your entire network.

1

PART

CoS Overview

[Overview](#) | 2

[Basic Concepts](#) | 15

[CoS Upgrade Requirements and Feature Change Compatibility](#) | 49

CHAPTER 1

Overview

IN THIS CHAPTER

- [Overview of Junos OS CoS | 2](#)
- [Overview of Policers | 5](#)

Overview of Junos OS CoS

IN THIS SECTION

- [CoS Standards | 3](#)
- [How Junos OS CoS Works | 4](#)
- [Default CoS Behavior | 5](#)

When a network experiences congestion and delay, some packets must be dropped. Junos OS *class of service* (CoS) enables you to divide traffic into classes and set various levels of throughput and packet loss when congestion occurs. You have greater control over packet loss because you can configure rules tailored to your needs.

You can configure CoS features to provide multiple classes of service for different applications. CoS also allows you to rewrite the Differentiated Services code point (DSCP) or IEEE 802.1p code-point bits of packets leaving an interface, thus allowing you to tailor packets for the network requirements of the remote peers.

CoS provides multiple classes of service for different applications. You can configure multiple forwarding classes for transmitting packets, define which packets are placed into each output queue, schedule the transmission service level for each queue, and manage congestion using a weighted random early detection (WRED) algorithm.

In designing CoS applications, you must carefully consider your service needs, and you must thoroughly plan and design your CoS configuration to ensure consistency and interoperability across all platforms in a CoS domain.

Because CoS is implemented in hardware rather than in software, you can experiment with and deploy CoS features without affecting packet forwarding and switching performance.

NOTE: CoS policies can be enabled or disabled on each switch interface. Also, each physical and *logical interface* on the switch can have associated custom CoS rules.

When you change or when you deactivate and then reactivate the class-of-service configuration, the system experiences packet drops because the system momentarily blocks traffic to change the mapping of incoming traffic to input queues.

This topic describes:

CoS Standards

The following RFCs define the standards for CoS capabilities:

- RFC 2474, *Definition of the Differentiated Services Field in the IPv4 and IPv6 Headers*
- RFC 2597, *Assured Forwarding PHB Group*
- RFC 2598, *An Expedited Forwarding PHB*
- RFC 2698, *A Two Rate Three Color Marker*
- RFC 3168, *The Addition of Explicit Congestion Notification (ECN) to IP*

The following data center bridging (DCB) standards are also supported to provide the CoS (and other characteristics) that Fibre Channel over Ethernet (FCoE) requires for transmitting storage traffic over an Ethernet network:

- IEEE 802.1Qbb, *priority-based flow control* (PFC)
- IEEE 802.1Qaz, enhanced transmission selection (ETS)
- IEEE 802.1AB (LLDP) extension called Data Center Bridging Capability Exchange Protocol (DCBX)

NOTE: OCX Series switches and NFX250 Network Services platforms do not support PFC and DCBX.

Juniper Networks QFX10000 switches support both enhanced transmission selection (ETS) hierarchical port scheduling and direct port scheduling.

How Junos OS CoS Works

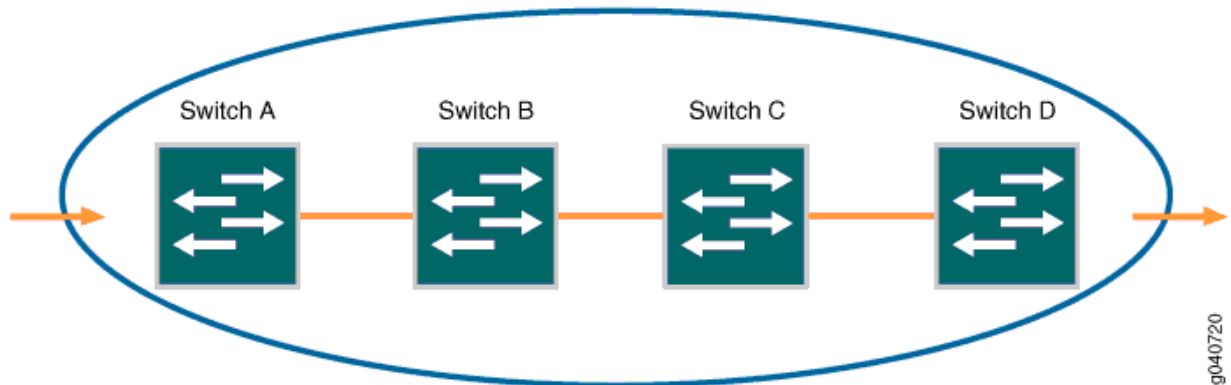
Junos OS CoS works by examining traffic entering the edge of your network. The switch classifies traffic into defined service groups to provide the special treatment of traffic across the network. For example, you can send voice traffic across certain links and data traffic across other links. In addition, the data traffic streams can be serviced differently along the network path to ensure that higher-paying customers receive better service. As the traffic leaves the network at the far edge, you can reclassify the traffic to meet the policies of the targeted peer by rewriting the DSCP or IEEE 802.1 code-point bits.

To support CoS, you must configure each switch in the network. Generally, each switch examines the packets that enter it to determine their CoS settings. These settings dictate which packets are transmitted first to the next downstream switch. Switches at the edges of the network might be required to alter the CoS settings of the packets that enter the network to classify the packets into the appropriate service groups.

In [Figure 1 on page 5](#), Switch A is receiving traffic. As each packet enters, Switch A examines the packet's current CoS settings and classifies the traffic into one of the groupings defined on the switch. This definition allows Switch A to prioritize its resources for servicing the traffic streams it receives. Switch A might alter the CoS settings (forwarding class and loss priority) of the packets to better match the defined traffic groups.

When Switch B receives the packets, it examines the CoS settings, determines the appropriate traffic groups, and processes the packet according to those settings. It then transmits the packets to Switch C, which performs the same actions. Switch D also examines the packets and determines the appropriate groups. Because Switch D sits at the far end of the network, it can reclassify (rewrite) the CoS code-point bits of the packets before transmitting them.

Figure 1: Packet Flow Across the Network



Default CoS Behavior

If you do not configure CoS settings, the software performs some CoS functions to ensure that the system forwards traffic and protocol packets with minimum delay when the network is experiencing congestion. Some CoS settings, such as classifiers, are automatically applied to each logical interface that you configure. Other settings, such as *rewrite rules*, are applied only if you explicitly associate them with an interface.

RELATED DOCUMENTATION

Overview of Policers

Understanding Junos CoS Components

Understanding CoS Packet Flow

Understanding CoS Hierarchical Port Scheduling (ETS)

Overview of Policers

IN THIS SECTION

- [Policer Overview | 6](#)
- [Policer Types | 9](#)
- [Policer Actions | 10](#)

- [Policer Colors | 11](#)
- [Filter-Specific Policers | 11](#)
- [Suggested Naming Convention for Policers | 11](#)
- [Policer Counters | 12](#)
- [Policer Algorithms | 12](#)
- [How Many Policers Are Supported? | 12](#)
- [Policers Can Limit Egress Firewall Filters | 13](#)

A switch polices traffic by limiting the input or output transmission rate of a class of traffic according to user-defined criteria. Policing (or rate-limiting) traffic allows you to control the maximum rate of traffic sent or received on an interface and to provide multiple priority levels or classes of service.

Policing is also an important component of firewall filters. You can achieve policing by including policers in *firewall filter* configurations.

Policer Overview

You use policers to apply limits to traffic flow and set consequences for packets that exceed these limits—usually applying a higher loss priority—so that if packets encounter downstream congestion, they can be discarded first. Policers apply only to unicast packets.

Policers provide two functions: metering and marking. A policer meters (measures) each packet against traffic rates and burst sizes that you configure. It then passes the packet and the metering result to the marker, which assigns a packet loss priority that corresponds to the metering result. [Figure 2 on page 8](#) illustrates this process.

NOTE: A policer restricts traffic at the configured transmission rate per PFE. In QFX10016, QFX10002, QFX10002-60C, and QFX10008 switches, when aggregated ethernet (AE) interface bundles span multiple PFEs, the overall transmission rate of the policer for the subscriber could exceed the configured transmission rate of the policer (depending on the number of PFEs involved).

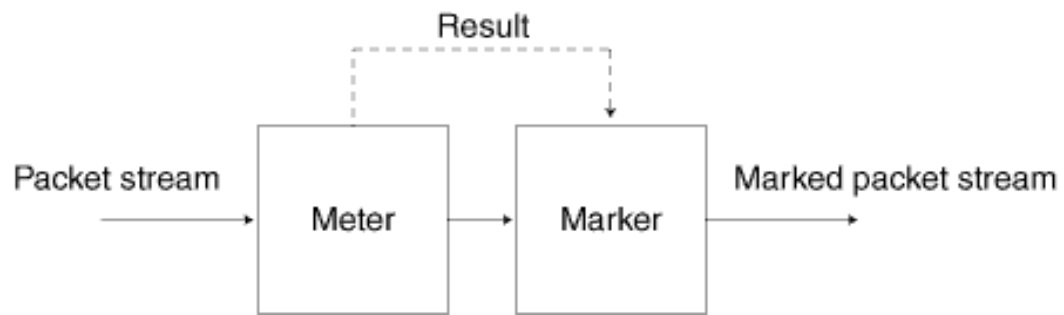
As an example:

- Policer with bandwidth-limit 100 mbps configured on an AE interface that has member links xe-1/0/0 (fpc1-pfe0) and xe-1/0/30 (fpc1-pfe1) . Here, the two member links belong to

FPC1, but are on different PFEs. When the policer is applied to the AE interface, this will result in a total bandwidth of 200 Mbps as policer is configured for two PFEs.

- Policer with bandwidth-limit 100 mbps configured on an AE interface that has member links xe-1/0/0 (fpc1-pfe0), et-2/0/1 (fpc2-pfe1) and xe-2/0/18:0 (fpc2-pfe2) . Here, one member link belongs to FPC1 and PFE0 on this FPC. The rest two member links belong to FPC2, but different PFEs. When the policer is applied to the AE interface, this will result in a total bandwidth of 300 Mbps as policer is configured for three PFEs.
- Policer with bandwidth-limit 100 mbps configured on an AE interface that has member links xe-1/0/0 and xe-1/0/1 on a single PFE (fpc1-pfe0) . Here, the member links belong to FPC1 and to the same PFE. When the policer is applied to the AE interface, this will result in a total bandwidth of 100 Mbps as policer is configured on a per PFE basis.

Figure 2: Flow of Tricolor Marking Policer Operation



g017049

After you name and configure a policer, you can use it by specifying it as an action in one or more firewall filters.

Policer Types

A switch supports three types of policers:

- **Single-rate two-color marker**—A two-color policer (or “policer” when used without qualification) meters the traffic stream and classifies packets into two categories of packet loss priority (PLP) according to a configured bandwidth and burst-size limit. You can mark packets that exceed the bandwidth and burst-size limit with a specified PLP or simply discard them.

You can specify this type of policer in an ingress or egress firewall.

NOTE: A two-color policer is most useful for metering traffic at the port (physical interface) level.

- **Single-rate three-color marker**—This type of policer is defined in RFC 2697, *A Single Rate Three Color Marker*, as part of an assured forwarding (AF) per-hop-behavior (PHB) classification system for a Differentiated Services (DiffServ) environment. This type of policer meters traffic based on one rate—the configured committed information rate (CIR) as well as the committed burst size (CBS) and the excess burst size (EBS). The CIR specifies the average rate at which bits are admitted to the switch. The CBS specifies the usual burst size in bytes and the EBS specifies the maximum burst size in bytes. The EBS must be greater than or equal to the CBS, and neither can be 0.

You can specify this type of policer in an ingress or egress firewall.

NOTE: A single-rate three-color marker (TCM) is most useful when a service is structured according to packet length and not peak arrival rate.

- **Two-rate three-color marker**—This type of policer is defined in RFC 2698, *A Two Rate Three Color Marker*, as part of an assured forwarding per-hop-behavior classification system for a Differentiated Services environment. This type of policer meters traffic based on two rates—the CIR and peak information rate (PIR) along with their associated burst sizes, the CBS and peak burst size (PBS). The PIR specifies the maximum rate at which bits are admitted to the network and must be greater than or equal to the CIR.

You can specify this type of policer in an ingress or egress firewall.

NOTE: A two-rate three-color policer is most useful when a service is structured according to arrival rates and not necessarily packet length.

See [Table 1 on page 10](#) for information about how metering results are applied for each of these policer types.

Policer Actions

Policer actions are implicit or explicit and vary by policer type. *Implicit* means that Junos OS assigns the loss priority automatically. [Table 1 on page 10](#) describes the policer actions.

Table 1: Policer Actions

Policer	Marking	Implicit Action	Configurable Action
Single-rate two-color	Green (conforming)	Assign low loss priority	None
	Red (nonconforming)	None	Discard
Single-rate three-color	Green (conforming)	Assign low loss priority	None
	Yellow (above the CIR and CBS)	Assign medium-high loss priority	None
	Red (above the EBS)	Assign high loss priority	Discard
Two-rate three-color	Green (conforming)	Assign low loss priority	None
	Yellow (above the CIR and CBS)	Assign medium-high loss priority	None
	Red (above the PIR and PBS)	Assign high loss priority	Discard

NOTE: If you specify a policer in an egress *firewall filter*, the only supported action is discard.

Policer Colors

Single-rate and two-rate three-color policers can operate in two modes:

- **Color-blind**—In color-blind mode, the three-color policer assumes that all packets examined have not been previously marked or metered. In other words, the three-color policer is “blind” to any previous coloring a packet might have had.
- **Color-aware**—In color-aware mode, the three-color policer assumes that all packets examined have been previously marked or metered. In other words, the three-color policer is “aware” of the previous coloring a packet might have had. In color-aware mode, the three-color policer can increase the PLP of a packet but cannot decrease it. For example, if a color-aware three-color policer meters a packet with a medium PLP marking, it can raise the PLP level to high but cannot reduce the PLP level to low.

Filter-Specific Policers

You can configure policers to be filter-specific, which means that Junos OS creates only one policer instance regardless of how many times the policer is referenced. When you do this on some QFX switches, rate limiting is applied in aggregate, so if you configure a policer to discard traffic that exceeds 1 Gbps and reference that policer in three different terms, the total bandwidth allowed by the filter is 1 Gbps. However, the behavior of a filter-specific policer is affected by how the firewall filter terms that reference the policer are stored in TCAM. If you create a filter-specific policer and reference it in multiple firewall filter terms, the policer allows more traffic than expected if the terms are stored in different TCAM slices. For example, if you configure a policer to discard traffic that exceeds 1 Gbps and reference that policer in three different terms that are stored in three separate memory slices, the total bandwidth allowed by the filter is 3 Gbps, not 1 Gbps. (This behavior does not occur in QFX10000 switches.)

To prevent this unexpected behavior from occurring, use the information about TCAM slices presented in [Planning the Number of Firewall Filters to Create](#) to organize your configuration file so that all the firewall filter terms that reference a given filter-specific policer are stored in the same TCAM slice.

Suggested Naming Convention for Policers

We recommend that you use the naming convention *policertypeTCM#-color type* when configuring three-color policers and *policer#* when configuring two-color policers. TCM stands for three-color marker. Because policers can be numerous and must be applied correctly to work, a simple naming convention makes it easier to apply the policers properly. For example, the first single-rate, color-aware three-color

policer configured would be named `srTCM1-ca`. The second two-rate, color-blind three-color configured would be named `trTCM2-cb`. The elements of this naming convention are explained below:

- `sr` (single-rate)
- `tr` (two-rate)
- `TCM` (tricolor marking)
- `1` or `2` (number of marker)
- `ca` (color-aware)
- `cb` (color-blind)

Policer Counters

On some QFX switches, each policer that you configure includes an implicit counter that counts the number of packets that exceed the rate limits that are specified for the policer. If you use the same policer in multiple terms—either within the same filter or in different filters—the implicit counter counts all the packets that are policed in all of these terms and provides the total amount. (This does not apply to QFX10000 switches.) If you want to obtain separate packet counts for each term on an affected switch, use these options:

- Configure a unique policer for each term.
- Configure only one policer, but use a unique, explicit counter in each term.

Policer Algorithms

Policing uses the *token-bucket algorithm*, which enforces a limit on average bandwidth while allowing bursts up to a specified maximum value. It offers more flexibility than the *leaky bucket algorithm* in allowing a certain amount of bursty traffic before it starts discarding packets.

NOTE: In an environment of light bursty traffic, QFX5200 might not replicate all multicast packets to two or more downstream interfaces. This occurs only at a line rate burst—if traffic is consistent, the issue does not occur. In addition, the issue occurs only when packet size increases beyond 6k in a one gigabit traffic flow.

How Many Policers Are Supported?

QFX10000 switches support 8K policers (all policer types). QFX5100 and QFX5200 switches support 1535 ingress policers and 1024 egress policers (assuming one policer per firewall filter term). QFX5110

switches support 6144 ingress policers and 1024 egress policers (assuming one policer per firewall filter term).

QFX3500 and QFX3600 standalone switches and QFabric Node devices support the following numbers of policers (assuming one policer per firewall filter term):

- Two-color policers used in ingress firewall filters: 767
- Three-color policers used in ingress firewall filters: 767
- Two-color policers used in egress firewall filters: 1022
- Three-color policers used in egress firewall filters: 512

Policers Can Limit Egress Firewall Filters

On some switches, the number of egress policers you configure can affect the total number of allowed egress firewall filters. Every policer has two implicit counters that take up two entries in a 1024-entry TCAM. These are used for counters, including counters that are configured as action modifiers in firewall filter terms. (Policers consume two entries because one is used for green packets and one is used for nongreen packets regardless of policer type.) If the TCAM becomes full, you are unable to commit any more egress firewall filters that have terms with counters. For example, if you configure and commit 512 egress policers (two-color, three-color, or a combination of both policer types), all of the memory entries for counters get used up. If later in your configuration file you insert additional egress firewall filters with terms that also include counters, *none* of the terms in those filters are committed because there is no available memory space for the counters.

Here are some additional examples:

- Assume that you configure egress filters that include a total of 512 policers and no counters. Later in your configuration file you include another egress filter with 10 terms, 1 of which has a counter action modifier. None of the terms in this filter are committed because there is not enough TCAM space for the counter.
- Assume that you configure egress filters that include a total of 500 policers, so 1000 TCAM entries are occupied. Later in your configuration file you include the following two egress filters:
 - Filter A with 20 terms and 20 counters. All the terms in this filter are committed because there is enough TCAM space for all the counters.
 - Filter B comes after Filter A and has five terms and five counters. *None* of the terms in this filter are committed because there is not enough memory space for *all* the counters. (Five TCAM entries are required but only four are available.)

You can prevent this problem by ensuring that egress firewall filter terms with counter actions are placed earlier in your configuration file than terms that include policers. In this circumstance, Junos OS commits

policers even if there is not enough TCAM space for the implicit counters. For example, assume the following:

- You have 1024 egress firewall filter terms with counter actions.
- Later in your configuration file you have an egress filter with 10 terms. None of the terms have counters but one has a policer action modifier.

You can successfully commit the filter with 10 terms even though there is not enough TCAM space for the implicit counters of the policer. The policer is committed without the counters.

RELATED DOCUMENTATION

Understanding Color-Blind Mode for Single-Rate Tricolor Marking

Understanding Color-Blind Mode for Two-Rate Tricolor Marking

Understanding Color-Aware Mode for Single-Rate Tricolor Marking

Understanding Color-Aware Mode for Two-Rate Tricolor Marking

Configuring Two-Color and Three-Color Policers to Control Traffic Rates

Basic Concepts

IN THIS CHAPTER

- [Configuring CoS | 15](#)
- [Understanding Junos CoS Components | 22](#)
- [Assigning CoS Components to Interfaces | 27](#)
- [Understanding CoS Packet Flow | 29](#)
- [Understanding Default CoS Settings | 33](#)
- [CoS Inputs and Outputs Overview | 47](#)

Configuring CoS

The traffic management class-of-service topics describe how to configure the Junos OS class-of-service (CoS) components. Junos CoS provides a flexible set of tools that enable you to fine tune control over the traffic on your network.

- Define classifiers that classify incoming traffic into forwarding classes to place traffic in groups for transmission.
- Map forwarding classes to output queues to define the type of traffic on each output queue.
- Configure schedulers for each output queue to control the service level (priority, bandwidth characteristics) of each type of traffic.
- Provide different service levels for the same forwarding classes on different interfaces.
- On switches that support data center bridging standards, configure lossless transport across the Ethernet network using priority-based flow control (PFC), Data Center Bridging Exchange protocol (DCBX), and enhanced transmission selection (ETS) hierarchical scheduling (OCX Series switches and NFX250 Network Services platform do not support lossless transport, PFC, and DCBX).
- Configure various CoS components individually or in combination to define CoS services.

NOTE: When you change the CoS configuration or when you deactivate and then reactivate the CoS configuration, the system experiences packet drops because the system momentarily blocks traffic to change the mapping of incoming traffic to input queues.

[Table 2 on page 17](#) lists the primary CoS configuration tasks by platform and provides links to those tasks.

NOTE: Links to features that are not supported on the platform for which you are looking up information might not be functional.

Table 2: CoS Configuration Tasks

CoS Configuration Task	Platforms Supported	Links
<p>Basic CoS Configuration:</p> <ul style="list-style-type: none"> Configure code-point aliases to assign a name to a pattern of code-point bits that you can use instead of the bit pattern when you configure CoS components such as classifiers and rewrite rules Configure classifiers and multidestination classifiers <ul style="list-style-type: none"> Set the forwarding class and loss priority of a packet based on the incoming CoS value and assign packets to output queues based on the associated forwarding class Change the host default output queue and mapping of DSCP bits used in the type of service (ToS) field Configure forwarding classes Configure rewrite rules to alter code point bit values in outgoing packets on the outbound interfaces of a switch so that the CoS treatment matches the policies of a targeted peer Configure Ethernet PAUSE flow control, a congestion relief feature that provides link-level flow control for all traffic on a full-duplex Ethernet link, including those that belong to Ethernet link aggregated (LAG) interfaces. On any particular interface, symmetric and asymmetric flow control are mutually exclusive. Assign the following CoS components to physical or logical interfaces: 	<ul style="list-style-type: none"> QFX3500 QFX3600 EX4600 NFX250 QFX5100 QFX5200 QFX5210 QFX10000 OCX1100 switches QFabric systems 	<ul style="list-style-type: none"> Defining CoS Code-Point Aliases (QFX10000 only) Example: Configuring Classifiers (Except QFX10000) Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p) (Except NFX250 and QFX10000) Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers Changing the Host Outbound Traffic Default Queue Mapping Example: Configuring Forwarding Classes Defining CoS Rewrite Rules (Except NFX250) Enabling and Disabling CoS Symmetric Ethernet PAUSE Flow Control (Except NFX250 and OCX1100) Configuring CoS Asymmetric Ethernet PAUSE Flow Control Assigning CoS Components to Interfaces

Table 2: CoS Configuration Tasks *(Continued)*

CoS Configuration Task	Platforms Supported	Links
<ul style="list-style-type: none"> Classifiers Congestion notification profiles Forwarding classes Forwarding class sets Output traffic control profiles Port schedulers Rewrite rules 		
<p>Configure Weighted random early detection (WRED) drop profiles that define the drop probability of packets of different packet loss probabilities (PLPs) as the output queue fills:</p> <ul style="list-style-type: none"> Configure WRED drop profiles where you associate WRED drop profiles with loss priorities in a scheduler. When you map the scheduler to a forwarding class (queue), you apply the interpolated drop profile to traffic of the specified loss priority on that queue. Configure drop profile maps that map a drop profile to a packet loss priority, and associate the drop profile and packet loss priority with a scheduler Configure explicit congestion notification (ECN) to enable end-to-end congestion notification between two endpoints on TCP/IP based networks. Apply WRED drop profiles to forwarding classes to control how the switch marks ECN-capable packets. 	<ul style="list-style-type: none"> QFX3500 QFX3600 EX4600 QFX5100 QFX5200 QFX5210 QFX10000 OCX1100 switches QFabric systems 	<ul style="list-style-type: none"> Example: Configuring WRED Drop Profiles Example: Configuring Drop Profile Maps Example: Configuring ECN

Table 2: CoS Configuration Tasks *(Continued)*

CoS Configuration Task	Platforms Supported	Links
<p>Configure queue schedulers and the bandwidth scheduling priority of individual queues. Schedulers define the CoS properties of output queues (output queues are mapped to forwarding classes, and classifiers map traffic into forwarding classes based on IEEE 802.1p or DSCP code points). Queue scheduling works with priority group scheduling to create a two-tier hierarchical scheduler. CoS scheduling properties include the amount of interface bandwidth assigned to the queue, the priority of the queue, whether explicit congestion notification (ECN) is enabled on the queue, and the WRED packet drop profiles associated with the queue.</p>	<ul style="list-style-type: none"> • QFX3500 • QFX3600 • EX4600 • NFX250 • QFX5100 • QFX5200 • QFX5210 • QFX10000 • OCX1100 switches • QFabric systems 	<ul style="list-style-type: none"> • (Except QFX10000) Example: Configuring Queue Schedulers • Example: Configuring Queue Scheduling Priority • (QFX10000 only) Example: Configuring Queue Schedulers for Port Scheduling
<p>Configure traffic control profiles to define the output bandwidth and scheduling characteristics of forwarding class sets (priority groups). The forwarding classes (queues) mapped to a forwarding class set share the bandwidth resources that you configure in the traffic control profile.</p>	<ul style="list-style-type: none"> • QFX3500 • QFX3600 • EX4600 • NFX250 • QFX5100 • QFX5200 • QFX5210 • QFX10000 • OCX1100 switches • QFabric systems 	<ul style="list-style-type: none"> • (Except NFX250) Defining CoS Traffic Control Profiles (Priority Group Scheduling) • (Except NFX250) Example: Configuring Traffic Control Profiles (Priority Group Scheduling) • Example: Configuring Minimum Guaranteed Output Bandwidth • (Except NFX250) Example: Configuring Maximum Output Bandwidth

Table 2: CoS Configuration Tasks *(Continued)*

CoS Configuration Task	Platforms Supported	Links
<p>Configure enhanced transmission selection (ETS) and forwarding class sets, and disable the ETS recommendation TLV. Hierarchical port scheduling, the Junos OS implementation of ETS, enables you to group priorities that require similar CoS treatment into priority groups. You define the port bandwidth resources for a priority group, and you define the amount of the priority group's resources that each priority in the group can use.</p>	<ul style="list-style-type: none"> • QFX3500 • QFX3600 • EX4600 • QFX5100 • OCX1100 switches • QFX10000 • QFabric systems 	<ul style="list-style-type: none"> • Example: Configuring Forwarding Class Sets • Example: Configuring CoS Hierarchical Port Scheduling (ETS) • (Except OCX1100) Disabling the ETS Recommendation TLV
<p>Configure Data Center Bridging Capability Exchange protocol (DCBX), which discovers the data center bridging (DCB) capabilities of peers by exchanging feature configuration information and is an extension of the Link Layer Discovery Protocol (LLDP)</p> <ul style="list-style-type: none"> • Configure the DCBX mode that an interface uses to communicate with the connected peer • Configure DCBX autonegotiation on a per-interface basis for each supported feature or application • Define each application for which you want DCBX to exchange application protocol information • Map applications to IEEE 802.1p code points • Apply an application map to a DCBX interface 	<ul style="list-style-type: none"> • QFX3500 • QFX3600 • EX4600 • QFX5100 • QFX5200 • QFX5210 • QFX10000 • QFabric systems 	<ul style="list-style-type: none"> • Example: Configuring DCBX Application Protocol TLV Exchange • Configuring the DCBX Mode • Configuring DCBX Autonegotiation • Defining an Application for DCBX Application Protocol TLV Exchange • Configuring an Application Map for DCBX Application Protocol TLV Exchange • Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange

Table 2: CoS Configuration Tasks *(Continued)*

CoS Configuration Task	Platforms Supported	Links
<p>Configure CoS for FCoE:</p> <ul style="list-style-type: none"> • Configure priority-based flow control (PFC) to divide traffic on one physical link into eight priorities • Configure a congestion notification profile (CNP) that enables priority-based flow control (PFC) on specified IEEE 802.1p priorities • Configure Multichassis link aggregation groups (MC-LAGs) to provide redundancy and load balancing between two switches • Configure two or more lossless forwarding classes and map them to different priorities • Configure lossless FCoE transport if your network uses a different priority than 3 • Configure multiple lossless FCoE priorities on a converged Ethernet network • If the FCoE network uses a different priority than priority 3 for FCoE traffic, configure a rewrite value to remap incoming traffic from the FC SAN to that priority after the interface encapsulates the FC packets in Ethernet • Configure lossless priorities for multiple types of traffic, such as FCoE and iSCSI 	<ul style="list-style-type: none"> • QFX3500 • QFX3600 • EX4600 • QFX5100 • QFX5200 • QFX5210 • QFX10000 • QFabric systems 	<ul style="list-style-type: none"> • Example: Configuring CoS PFC for FCoE Traffic • Example: Configuring CoS for FCoE Transit Switch Traffic Across an MC-LAG • Configuring CoS PFC (Congestion Notification Profiles) • (QFX3500 and QFabric only) "Example: Configuring IEEE 802.1p Priority Remapping on an FCoE-FC Gateway" on page 624 • Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces • Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch) • Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface • (QFX3500, NFX250, and QFabric only) "Configuring CoS Fixed Classifier Rewrite Values for Native FC Interfaces (NP_Ports)" on page 156 • Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI)

Understanding Junos CoS Components

IN THIS SECTION

- [Code-Point Aliases | 22](#)
- [Policers | 22](#)
- [Classifiers | 22](#)
- [Forwarding Classes | 23](#)
- [Forwarding Class Sets | 24](#)
- [Flow Control \(Ethernet PAUSE, PFC, and ECN\) | 25](#)
- [WRED Profiles and Tail Drop | 26](#)
- [Schedulers | 26](#)
- [Rewrite Rules | 27](#)

This topic describes the Junos OS class-of-service (CoS) components:

Code-Point Aliases

A *code-point alias* assigns a name to a pattern of code-point bits. You can use this name instead of the bit pattern when you configure other CoS components such as classifiers and *rewrite rules*.

Policers

Policers limit traffic of a certain class to a specified bandwidth and burst size. Packets exceeding the policer limits can be discarded, or can be assigned to a different forwarding class, a different loss priority, or both. You define policers with filters that you can associate with input interfaces.

Classifiers

Packet classification associates incoming packets with a particular CoS servicing level. In Junos OS, *classifiers* associate packets with a forwarding class and loss priority and assign packets to output queues based on the associated forwarding class. Junos OS supports two general types of classifiers:

- Behavior aggregate (BA) or CoS value traffic classifiers—Examine the CoS value in the packet header. The value in this single field determines the CoS settings applied to the packet. BA classifiers allow

you to set the forwarding class and loss priority of a packet based on the Differentiated Services code point (DSCP) value, IEEE 802.1p value, or MPLS EXP value.

NOTE: OCX Series switches and NFX250 Network Services platform do not support MPLS.

- **Multifield traffic classifiers**—Examine multiple fields in the packet, such as source and destination addresses and source and destination port numbers of the packet. With multifield classifiers, you set the forwarding class and loss priority of a packet based on *firewall filter* rules.

On switches that require the separation of unicast and multideestination (multicast, broadcast, and destination lookup fail) traffic, you create separate unicast classifiers and multideestination classifiers. You cannot assign unicast traffic and multideestination traffic to the same classifier. You can apply unicast classifiers to one or more interfaces. Multideestination classifiers apply to all of the switch interfaces and cannot be applied to individual interfaces. Switches that require the separation of unicast and multideestination traffic have 12 output queues to provide 4 output queues reserved for multideestination traffic.

On switches that do not separate unicast and multideestination traffic, unicast and multideestination traffic use the same classifiers, and you do not create a separate special classifier for multideestination traffic. Switches that do not separate unicast and multideestination traffic have eight output queues because no extra queues are required to separate the traffic.

Forwarding Classes

Forwarding classes group packets for transmission and CoS. You assign each packet to an output queue based on the packet's forwarding class. Forwarding classes affect the forwarding, scheduling, and rewrite marking policies applied to packets as they transit the switch.

Switches provide up to five default forwarding classes:

- **best-effort**—Best-effort traffic
- **fcoe**—Fibre Channel over Ethernet traffic
- **no-loss**—Lossless traffic
- **network-control**—Network control traffic
- **mcast**—Multicast traffic

NOTE: The default `mcast` forwarding class applies only to switches that require the separation of unicast and multideestination (multicast, broadcast, and destination lookup fail) traffic. On these

switches, you create separate forwarding classes for the two types of traffic. The default mcast forwarding class transports only multidestination traffic, and the default best-effort, fcoe, no-loss, and network-control forwarding classes transport only unicast traffic. Unicast forwarding classes map to unicast output queues, and multidestination forwarding classes map to multidestination output queues. You cannot assign unicast traffic and multidestination traffic to the same forwarding class or to the same output queue. Switches that require the separation of unicast and multidestination traffic have 12 output queues, 8 for unicast traffic and 4 for multidestination traffic.

On switches that do not separate unicast and multidestination traffic, unicast and multidestination traffic use the same forwarding classes and output queues, so the mcast forwarding class is not valid. You do not create separate forwarding classes for multidestination traffic. Switches that do not separate unicast and multidestination traffic have eight output queues because no extra queues are required to separate the traffic.

NOTE: On OCX Series switches only, do not map traffic to the default fcoe and no-loss forwarding classes. By default, the DSCP default classifier does not map traffic to the fcoe and no-loss forwarding classes, so by default, OCX Series switches do not classify traffic into those forwarding classes. (On other switches, the fcoe and no-loss forwarding classes provide lossless transport for Layer 2 traffic. OCX Series switches do not support lossless Layer 2 transport.)

Switches support a total of either 12 forwarding classes (8 unicast forwarding classes and 4 multicast forwarding classes), or 8 forwarding classes (unicast and multidestination traffic use the same forwarding classes), which provides flexibility in classifying traffic.

NFX250 Network Services platform provide the following forwarding classes:

- best-effort (be)—Provides no service profile. Loss priority is typically not carried in a CoS value.
- expedited-forwarding (ef)—Provides a low loss, low latency, low jitter, assured bandwidth, end-to-end service.
- assured-forwarding (af)—Provides a group of values you can define and includes four subclasses: AF1, AF2, AF3, and AF4, each with two drop probabilities: low and high.
- network-control (nc)—Supports protocol control and thus is typically high priority.

Forwarding Class Sets

You can group forwarding classes (output queues) into *forwarding class sets* to apply CoS to groups of traffic that require similar treatment. Forwarding class sets map traffic into priority groups to support enhanced transmission selection (ETS), which is described in IEEE 802.1Qaz.

You can configure up to three unicast forwarding class sets and one multicast forwarding class set. For example, you can configure different forwarding class sets to apply CoS to unicast groups of local area network (LAN) traffic, storage area network (SAN) traffic, and high-performance computing (HPC) traffic, and configure another group for multicast traffic.

Within each forwarding class set, you can configure special CoS treatment for the traffic mapped to each individual queue. This provides the ability to configure CoS in a two-tier hierarchical manner. At the forwarding class set tier, you configure CoS for groups of traffic using a *traffic control profile*. At the queue tier, you configure CoS for individual output queues within a forwarding class set using a *scheduler* that you map to a queue (forwarding class) using a *scheduler map*.

Flow Control (Ethernet PAUSE, PFC, and ECN)

Ethernet PAUSE (described in IEEE 802.3X) is a link-level flow control mechanism. During periods of network congestion, Ethernet PAUSE stops all traffic on a full-duplex Ethernet link for a period of time specified in the PAUSE message.

NOTE: QFX10000 switches do not support Ethernet PAUSE.

Priority-based flow control (PFC) is described in IEEE 802.1Qbb as part of the IEEE data center bridging (DCB) specifications for creating a lossless Ethernet environment to transport loss-sensitive flows such as Fibre Channel over Ethernet (FCoE) traffic.

NOTE: OCX Series switches do not support PFC.

PFC is a link-level flow control mechanism similar to Ethernet PAUSE. However, Ethernet PAUSE stops all traffic on a link for a period of time. PFC decouples the pause function from the physical link and divides the traffic on the link into eight priorities (3-bit IEEE 802.1p code points). You can think of the eight priorities as eight “lanes” of traffic. You can apply pause selectively to the traffic on any priority without pausing the traffic on other priorities on the same link.

The granularity that PFC provides allows you to configure different levels of CoS for different types of traffic on the link. You can create lossless lanes for traffic such as FCoE, LAN backup, or management, while using standard frame-drop methods of congestion management for IP traffic on the same link.

NOTE: If you transport FCoE traffic, you must enable PFC on the priority assigned to FCoE traffic (usually IEEE 802.1p code point 011 on interfaces that carry FCoE traffic).

Explicit congestion notification (ECN) enables end-to-end congestion notification between two endpoints on TCP/IP based networks. ECN must be enabled on both endpoints and on all of the intermediate devices between the endpoints for ECN to work properly. Any device in the transmission path that does not support ECN breaks the end-to-end ECN functionality. ECN notifies networks about congestion with the goal of reducing packet loss and delay by making the sending device decrease the transmission rate until the congestion clears, without dropping packets. RFC 3168, *The Addition of Explicit Congestion Notification (ECN) to IP*, defines ECN.

WRED Profiles and Tail Drop

A weighted random early detection (WRED) profile (drop profile) defines parameters that enable the network to drop packets during periods of congestion. A *drop profile* defines the conditions under which packets of different loss priorities drop, by determining the probability of dropping a packet for each loss priority when output queues become congested. Drop profiles essentially set a value for a level of queue fullness—when the queue fills to the level of the queue fullness value, packets drop. The combination of queue fill level, the probability of dropping a packet at that fill level, and loss priority of the packet, determine whether a packet is dropped or forwarded. Each pairing of a fill level with a drop probability creates a point on a drop profile curve.

You can associate different drop profiles with different loss priorities to set the probability of dropping packets. You can apply a drop profile for each loss priority to a forwarding class (output queue) by applying a drop profile to a scheduler, and then mapping the scheduler to a forwarding class using a scheduler map. When the queue mapped to the forwarding class experiences congestion, the drop profile determines the level of packet drop for traffic of each loss priority in that queue.

Loss priority affects the scheduling of a packet without affecting the packet's relative ordering. Typically you mark packets exceeding a particular service level with a high loss priority.

Tail drop is a simple drop mechanism that drops all packets indiscriminately during periods of congestion, without differentiating among the packet loss priorities of traffic flows. Tail drop requires only one curve point that corresponds to the maximum depth of the output queue, and drop probability when traffic exceeds the buffer depth is 100 percent (all packets that cannot be stored in the queue are dropped). WRED is superior to tail-drop because WRED enables you to treat traffic of different priorities in a differentiated manner, so that higher priority traffic receives preference, and because of the ability to set multiple points on the drop curve.

Schedulers

Each switch interface has multiple queues assigned to store packets. The switch determines which queue to service based on a particular method of scheduling. This process often involves determining the sequence in which different types of packets should be transmitted.

You can define the scheduling priority (priority), minimum guaranteed bandwidth (transmit-rate), maximum bandwidth (shaping-rate), and WRED profiles to be applied to a particular queue (forwarding

class) for packet transmission. By default, extra bandwidth is shared among queues in proportion to the minimum guaranteed bandwidth of each queue. On switches that support the `excess-rate` statement, you can configure the percentage of shared extra bandwidth an output queue receives independently from the minimum guaranteed bandwidth transmit rate, or you can use default bandwidth sharing based on the transmit rate.

A scheduler map associates a specified forwarding class with a scheduler configuration. You can associate up to four user-defined scheduler maps with the interfaces.

Rewrite Rules

A *rewrite rule* sets the appropriate CoS bits in the outgoing packet. This allows the next downstream device to classify the packet into the appropriate service group. Rewriting (marking) outbound packets is useful when the switch is at the border of a network and must change the CoS values to meet the policies of the targeted peer.

NOTE: Ingress firewall filters can also rewrite forwarding class and loss priority values.

RELATED DOCUMENTATION

| *Understanding CoS Packet Flow*

Assigning CoS Components to Interfaces

After you define the following CoS components, you assign them to physical or logical interfaces. Components that you assign to physical interfaces are valid for all of the logical interfaces configured on the physical interface. Components that you assign to a logical interface are valid only for that logical interface.

- Classifiers—Assign only to logical interfaces; on some switches, you can apply classifiers to physical Layer 3 interfaces and the classifiers are applied to all logical interfaces on the physical interface.
- Congestion notification profiles—Assign only to physical interfaces.

NOTE: OCX Series switches and NFX250 Network Services platform do not support congestion notification profiles.

- Forwarding classes—Assign to interfaces by mapping to forwarding class sets.
- Forwarding class sets—Assign only to physical interfaces.
- Output traffic control profiles—Assign only to physical interfaces (with a forwarding class set).
- Port schedulers—Assign only to physical interfaces on switches that support port scheduling. Associate the scheduler with a forwarding class in a scheduler map and apply the scheduler map to the physical interface.
- Rewrite rules—Assign only to logical interfaces; on some switches, you can apply classifiers to physical Layer 3 interfaces and the classifiers are applied to all logical interfaces on the physical interface.

You can assign a CoS component to a single interface or to multiple interfaces using wildcards. You can also assign a congestion notification profile or a forwarding class set globally to all interfaces.

To assign CoS components to interfaces:

Assign a CoS component to a physical interface by associating a CoS component (for example, a forwarding class set named `be-priority-group`) with an interface:

```
[edit class-of-service interfaces]
user@switch# set xe-0/0/7 forwarding-class-set be-priority-group
```

Assign a CoS component to a logical interface by associating a CoS component (for example, a classifier named `be_classifier`) with a logical interface:

```
[edit class-of-service interfaces]
user@switch# set xe-0/0/7 unit 0 classifiers dscp be_classifier
```

Assign a CoS component to multiple interfaces by associating a CoS component (for example, a rewrite rule named `customup-rw`) to all 10-Gigabit Ethernet interfaces on the switch, use wildcard characters for the interface name and logical interface (unit) number:

```
[edit class-of-service interfaces]
user@switch# set xe-* unit * rewrite-rules ieee-802.1 customup-rw xe-* unit * rewrite-rules
ieee-802.1 customup-rw
```

Assign a congestion notification profile or a forwarding class set globally to all interfaces using the `set class-of-service interfaces all` statement. For example, to assign a forwarding class set named `be-priority-group` to all interfaces:

```
[edit class-of-service interfaces]
user@switch# set all forwarding-class-set be-priority-group
```

NOTE: If there is an existing CoS configuration of any type on an interface, the global configuration is not applied to that particular interface. The global configuration is applied to all interfaces that do not have an existing CoS configuration.

For example, if you configure a rewrite rule, assign it to interfaces `xe-0/0/20.0` and `xe-0/0/22.0`, and then configure a forwarding class set and apply it to all interfaces, the forwarding class set is applied to every interface except `xe-0/0/20` and `xe-0/0/22`.

RELATED DOCUMENTATION

Monitoring Interfaces That Have CoS Components

Understanding Junos CoS Components

Understanding CoS Packet Flow

When a packet traverses a switch, the switch provides the appropriate level of service to the packet using either default *class-of-service* (CoS) settings or CoS settings that you configure. On ingress ports, the switch classifies packets into appropriate forwarding classes and assigns a loss priority to the packets. On egress ports, the switch applies packet scheduling and (if you have configured them) *rewrite rules* to re-mark packets.

You can configure CoS on Layer 2 logical interfaces, and you can configure CoS on Layer 3 physical interfaces if you have defined at least one *logical interface* on the Layer 3 physical interface. You cannot configure CoS on Layer 2 physical interfaces and Layer 3 logical interfaces.

For Layer 2 traffic, either use the default CoS settings or configure CoS on each logical interface. You can apply different CoS settings to different Layer 2 logical interfaces.

NOTE: OCX Series switches do not support Layer 2 interfaces (family ethernet-switching).

For Layer 3 traffic, either use the default CoS settings or configure CoS on the physical interface (not on the logical unit). The switch uses the CoS applied on the physical Layer 3 interface for all logical Layer 3 interfaces configured on the physical Layer 3 interface.

The switch applies CoS to packets as they flow through the system:

- An interface has one or more classifiers of different types applied to it (configure this at the [edit class-of-service interfaces] hierarchy level). The classifier types are based on the portion of the incoming packet that the classifier examines (IEEE 802.1p code point bits or DSCP code point bits).
- When a packet enters an ingress port, the classifier assigns the packet to a forwarding class and a loss priority based on the code point bits of the packet (configure this at the [edit class-of-service classifiers] hierarchy level).
- The switch assigns each forwarding class to an output queue (configure this at the [edit class-of-service forwarding-classes] hierarchy level).
- Input (and output) policers meter traffic and can change the forwarding class and loss priority if a traffic flow exceeds its service level.
- A scheduler map is applied to each interface. When a packet exits an egress port, the scheduler map controls how it is treated (configure this at the [edit class-of-service interfaces] hierarchy level). A scheduler map assigns schedulers to forwarding classes (configure this at the [edit class-of-service scheduler-maps] hierarchy level).
- A scheduler defines how traffic is treated at the egress interface output queue (configure this at the [edit class-of-service schedulers] hierarchy level). You control the transmit rate, shaping rate, priority, and drop profile of each forwarding class by mapping schedulers to forwarding classes in scheduler maps, then applying scheduler maps to interfaces.
- A drop-profile defines how aggressively to drop packets that are mapped to a particular scheduler (configure this at the [edit class-of-service drop-profiles] hierarchy level).
- A rewrite rule takes effect as the packet leaves an interface that has a rewrite rule configured (configure this at the [edit class-of-service rewrite-rules] hierarchy level). The rewrite rule writes information to the packet (for example, a rewrite rule can re-mark the code point bits of outgoing traffic) according to the forwarding class and loss priority of the packet.

Figure 3 on page 31 is a high-level flow diagram of how packets from various sources enter switch interfaces, are classified at the ingress, and then scheduled (provided bandwidth) at the egress queues.

Figure 3: CoS Classifier, Queues, and Scheduler

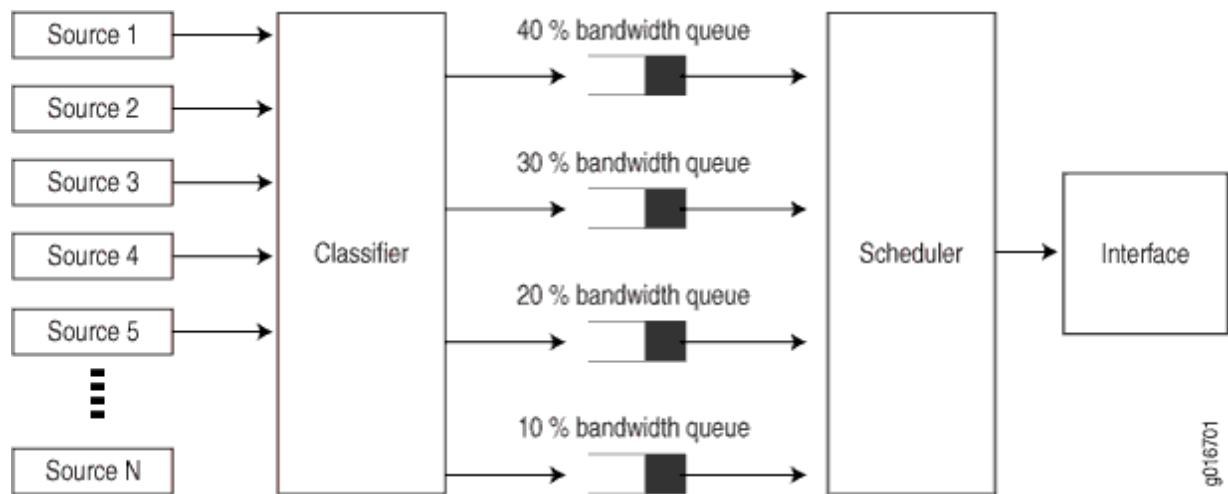
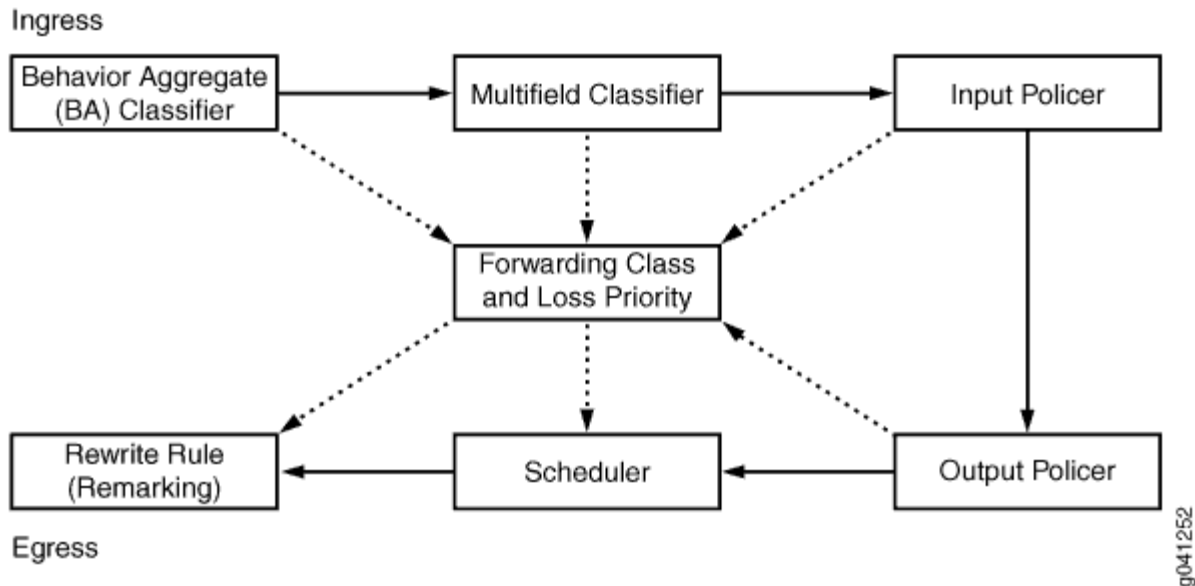


Figure 4 on page 32 shows the packet flow through the CoS components that you can configure.

Figure 4: Packet Flow Through Configurable CoS Components



The middle box (Forwarding Class and Loss Priority) represents two values that you can use on ingress and egress interfaces. The system uses these values for classifying traffic on ingress interfaces and for rewrite rule re-marking on egress interfaces. Each outer box represents a process component. The components in the top row apply to incoming packets. The components in the bottom row apply to outgoing packets.

The solid-line arrows show the direction of packet flow from ingress to egress. The dotted-line arrows that point to the forwarding class and loss priority box indicate processes that configure (set) the forwarding class and loss priority. The dotted-line arrows that point away from the forwarding class and loss priority box indicate processes that use forwarding class and loss priority as input values on which to base actions.

For example, the BA classifier sets the forwarding class and loss priority of incoming packets, so the forwarding class and loss priority are outputs of the classifier and the arrow points away from the classifier. The scheduler receives the forwarding class and loss priority settings, and queues the outgoing packets based on those settings, so the arrow points toward the scheduler.

Understanding Default CoS Settings

IN THIS SECTION

- [Default Forwarding Classes and Queue Mapping | 33](#)
- [Default Forwarding Class Sets \(Priority Groups\) | 34](#)
- [Default Code-Point Aliases | 35](#)
- [Default Classifiers | 37](#)
- [Default Rewrite Rules | 42](#)
- [Default Drop Profile | 42](#)
- [Default Schedulers | 42](#)
- [Default Scheduler Maps | 46](#)
- [Default Shared Buffer Configuration | 46](#)

If you do not configure CoS settings, Junos OS performs some CoS functions to ensure that traffic and protocol packets are forwarded with minimum delay when the network experiences congestion. Some default mappings are automatically applied to each *logical interface* that you configure.

You can display default CoS settings by issuing the `show class-of-service operational mode command`.

This topic describes the default configurations for the following CoS components:

Default Forwarding Classes and Queue Mapping

[Table 3 on page 33](#) shows the default mapping of the default forwarding classes to queues and packet drop attribute.

Table 3: Default Forwarding Classes and Queue Mapping

Default Forwarding Class	Description	Default Queue Mapping	Packet Drop Attribute
best-effort (be)	Best-effort traffic class (priority 0, IEEE 802.1p code point 000)	0	drop

Table 3: Default Forwarding Classes and Queue Mapping (Continued)

Default Forwarding Class	Description	Default Queue Mapping	Packet Drop Attribute
fcoe	Guaranteed delivery for FCoE traffic (priority 3, IEEE 802.1p code point 011)	3	no-loss
no-loss	Guaranteed delivery for TCP no-loss traffic (priority 4, IEEE 802.1p code point 100)	4	no-loss
network-control (nc)	Network control traffic (priority 7, IEEE 802.1p code point 111)	7	drop
(Excluding QFX10000) mcast	Multidestination traffic	8	drop NOTE: You cannot configure multidestination forwarding classes as no-loss (lossless) traffic classes.

NOTE: On the QFX10000 switch, unicast and multidestination (multicast, broadcast, and destination lookup fail) traffic use the same forwarding classes and output queues 0 through 7.

Default Forwarding Class Sets (Priority Groups)

If you do not explicitly configure forwarding class sets, the system automatically creates a default forwarding class set that contains all of the forwarding classes on the switch. The system assigns 100 percent of the port output bandwidth to the default forwarding class set.

Ingress traffic is classified based on the default classifier settings. The forwarding classes (queues) in the default forwarding class set receive bandwidth based on the default scheduler settings. Forwarding classes that are not part of the default scheduler receive no bandwidth.

The default forwarding class set is transparent. It does not appear in the configuration and is used for Data Center Bridging Capability Exchange (DCBX) protocol advertisement.

Default Code-Point Aliases

[Table 4 on page 35](#) shows the default mapping of code-point aliases to IEEE code points.

Table 4: Default IEEE 802.1 Code-Point Aliases

CoS Value Types	Mapping
be	000
be1	001
ef	010
ef1	011
af11	100
af12	101
nc1	110
nc2	111

[Table 5 on page 35](#) shows the default mapping of code-point aliases to DSCP and DSCP IPv6 code points.

Table 5: Default DSCP and DCSP IPv6 Code-Point Aliases

CoS Value Types	Mapping
ef	101110
af11	001010

Table 5: Default DSCP and DCSP IPv6 Code-Point Aliases *(Continued)*

CoS Value Types	Mapping
af12	001100
af13	001110
af21	010010
af22	010100
af23	010110
af31	011010
af32	011100
af33	011110
af41	100010
af42	100100
af43	100110
be	000000
cs1	001000
cs2	010000
cs3	011000

Table 5: Default DSCP and DCSP IPv6 Code-Point Aliases (Continued)

CoS Value Types	Mapping
cs4	100000
cs5	101000
nc1	110000
nc2	111000

Default Classifiers

The switch applies default unicast IEEE 802.1, unicast DSCP, and multidestination classifiers to each interface that does not have explicitly configured classifiers. If you explicitly configure one type of classifier but not other types of classifiers, the system uses only the configured classifier and does not use default classifiers for other types of traffic.

NOTE: The QFX10000 switch applies the default MPLS EXP classifier to a logical interface if you enable the MPLS protocol family on that interface.

There are two different default unicast IEEE 802.1 classifiers, a trusted classifier for ports that are in trunk mode or tagged-access mode, and an untrusted classifier for ports that are in access mode. [Table 6 on page 37](#) shows the default mapping of IEEE 802.1 code-point values to forwarding classes and loss priorities for ports in trunk mode or tagged-access mode.

Table 6: Default IEEE 802.1 Classifiers for Ports in Trunk Mode or Tagged Access Mode (Trusted Classifier)

Code Point	Forwarding Class	Loss Priority
be (000)	best-effort	low
be1 (001)	best-effort	low

Table 6: Default IEEE 802.1 Classifiers for Ports in Trunk Mode or Tagged Access Mode (Trusted Classifier) (Continued)

Code Point	Forwarding Class	Loss Priority
ef (010)	best-effort	low
ef1 (011)	fcoe	low
af11 (100)	no-loss	low
af12 (101)	best-effort	low
nc1 (110)	network-control	low
nc2 (111)	network-control	low

[Table 7 on page 38](#) shows the default mapping of IEEE 802.1p code-point values to forwarding classes and loss priorities for ports in access mode (all incoming traffic is mapped to best-effort forwarding classes).

Table 7: Default IEEE 802.1 Classifiers for Ports in Access Mode (Untrusted Classifier)

Code Point	Forwarding Class	Loss Priority
000	best-effort	low
001	best-effort	low
010	best-effort	low
011	best-effort	low
100	best-effort	low
101	best-effort	low

Table 7: Default IEEE 802.1 Classifiers for Ports in Access Mode (Untrusted Classifier) (Continued)

Code Point	Forwarding Class	Loss Priority
110	best-effort	low
111	best-effort	low

[Table 8 on page 39](#) shows the default mapping of IEEE 802.1 code-point values to multidestination (multicast, broadcast, and destination lookup fail traffic) forwarding classes and loss priorities.

Table 8: Default IEEE 802.1 Multidestination Classifiers

Code Point	Forwarding Class	Loss Priority
be (000)	mcast	low
be1 (001)	mcast	low
ef (010)	mcast	low
ef1 (011)	mcast	low
af11 (100)	mcast	low
af12 (101)	mcast	low
nc1 (110)	mcast	low
nc2 (111)	mcast	low

[Table 9 on page 40](#) shows the default mapping of DSCP code-point values to forwarding classes and loss priorities for DSCP IP and DCSP IPv6.

NOTE: There are no default DSCP IP classifiers for multideestination traffic. DSCP IPv6 classifiers are not supported for multideestination traffic.

Table 9: Default DSCP IP and IPv6 Classifiers

Code Point	Forwarding Class	Loss Priority
ef (101110)	best-effort	low
af11 (001010)	best-effort	low
af12 (001100)	best-effort	low
af13 (001110)	best-effort	low
af21 (010010)	best-effort	low
af22 (010100)	best-effort	low
af23 (010110)	best-effort	low
af31 (011010)	best-effort	low
af32 (011100)	best-effort	low
af33 (011110)	best-effort	low
af41 (100010)	best-effort	low
af42 (100100)	best-effort	low
af43 (100110)	best-effort	low

Table 9: Default DSCP IP and IPv6 Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
be (000000)	best-effort	low
cs1 (001000)	best-effort	low
cs2 (010000)	best-effort	low
cs3 (011000)	best-effort	low
cs4 (100000)	best-effort	low
cs5 (101000)	best-effort	low
nc1 (110000)	network-control	low
nc2 (111000)	network-control	low

On QFX10000 switches, [Table 10 on page 41](#) shows the default mapping of MPLS EXP code-point values to forwarding classes and loss priorities.

Table 10: Default EXP Classifiers on QFX10000 Switches

Code Point	Forwarding Class	Loss Priority
000	best-effort	low
001	best-effort	high
010	expedited-forwarding	low
011	expedited-forwarding	high
100	assured-forwarding	low

Table 10: Default EXP Classifiers on QFX10000 Switches (Continued)

Code Point	Forwarding Class	Loss Priority
101	assured-forwarding	high
110	network-control	low
111	network-control	high

Default Rewrite Rules

There are no default *rewrite rules*. If you do not explicitly configure rewrite rules, the switch does not reclassify egress traffic.

Default Drop Profile

[Table 11 on page 42](#) shows the default drop profile configuration.

Table 11: Default Drop Profile

Fill Level	Drop Probability
100	100

Default Schedulers

[Table 12 on page 43](#) shows the default scheduler configuration.

Table 12: Default Schedulers

Default Scheduler and Queue Number	Transmit Rate (Guaranteed Minimum Bandwidth)	Shaping Rate (Maximum Bandwidth)	Excess Bandwidth Sharing	Priority	Buffer Size
best-effort forwarding class scheduler (queue 0)	5% (QFX10000 15%)	None	5% (QFX10000 15%)	low	5% (QFX10000 15%)
fcoe forwarding class scheduler (queue 3)	35%	None	35%	low	35%
no-loss forwarding class scheduler (queue 4)	35%	None	35%	low	35%
network-control forwarding class scheduler (queue 7)	5% (QFX10000 15%)	None	5% (QFX10000 15%)	low	5% (QFX10000 15%)
(Excluding QFX10000) mcast forwarding class scheduler (queue 8)	20%	None	20%	low	20%

NOTE: The minimum guaranteed bandwidth (transmit rate) also determines the amount of excess (extra) bandwidth that the queue can share. Extra bandwidth is allocated to queues in proportion to the transmit rate of each queue. On QFX10000 switches, you can use the `excess-rate` statement to override the default transmit rate setting and configure the excess bandwidth percentage independently of the transmit rate.

By default, only the five default schedulers shown in [Table 12 on page 43](#), excluding the mcast scheduler on QFX10000 switches, have traffic mapped to them. Only the queues associated with the default schedulers, and forwarding classes on QFX10000 switches, receive default bandwidth, based on the default scheduler transmit rate. (You can configure schedulers and forwarding classes to allocate bandwidth to other queues or to change the default bandwidth of a default queue.) In addition, other than on QFX5200, QFX5210, and QFX10000 switches, multidestination queue 11 receives enough bandwidth from the default multidestination scheduler to handle CPU-generated multidestination

traffic. If a forwarding class does not transport traffic, the bandwidth allocated to that forwarding class is available to other forwarding classes.

NOTE: On QFX10000 switches, unicast and multdestination (multicast, broadcast, and destination lookup fail) traffic use the same forwarding classes and output queues.

Default hierarchical scheduling, known as enhanced transmission selection (ETS, defined in IEEE 802.1Qaz), divides the total port bandwidth between two groups of traffic: unicast traffic and multdestination traffic. By default, unicast traffic consists of queue 0 (best-effort forwarding class), queue 3 (fcoe forwarding class), queue 4 (no-loss forwarding class), and queue 7 (network-control forwarding class). Unicast traffic receives and shares a total of 80 percent of the port bandwidth. By default, multdestination traffic (mcast queue 8) receives a total of 20 percent of the port bandwidth. So on a 10-Gigabit port, default scheduling provides unicast traffic 8-Gbps of bandwidth and multdestination traffic 2-Gbps of bandwidth.

NOTE: Except on QFX5200, QFX5210, and QFX10000 switches, multdestination queue 11 also receives a small amount of default bandwidth from the multdestination scheduler. CPU-generated multdestination traffic uses queue 11, so you might see a small number of packets egress from queue 11. In addition, in the unlikely case that firewall filter match conditions map multdestination traffic to a unicast forwarding class, that traffic uses queue 11.

On QFX10000 switches, default scheduling is port scheduling. Default hierarchical scheduling, known as ETS, allocates the total port bandwidth to the four default forwarding classes served by the four default schedulers, as defined by the four default schedulers. The result is the same as direct port scheduling. Configuring hierarchical port scheduling, however, enables you to group forwarding classes that carry similar types of traffic into forwarding class sets (also called priority groups), and to assign port bandwidth to each forwarding class set. The port bandwidth assigned to the forwarding class set is then assigned to the forwarding classes within the forwarding class set. This hierarchy enables you to control port bandwidth allocation with greater granularity, and enables hierarchical sharing of extra bandwidth to better utilize link bandwidth.

Default scheduling for all switches uses weighted round-robin (WRR) scheduling. Each queue receives a portion (weight) of the total available interface bandwidth. The scheduling weight is based on the transmit rate of the default scheduler for that queue. For example, queue 7 receives a default scheduling weight of 5 percent, 15 percent on QFX10000 switches, of the available bandwidth, and queue 4 receives a default scheduling weight of 35 percent of the available bandwidth. Queues are mapped to forwarding classes (for example, queue 7 is mapped to the network-control forwarding class and queue 4 is mapped to the no-loss forwarding class), so forwarding classes receive the default bandwidth for the queues to which they are mapped. Unused bandwidth is shared with other default queues.

If you want non-default (unconfigured) queues to forward traffic, you should explicitly map traffic to those queues (configure the forwarding classes and queue mapping) and create schedulers to allocate bandwidth to those queues. For example, except on QFX5200, QFX5210, and QFX10000 switches, by default, queues 1, 2, 5, and 6 are unconfigured, and multdestination queues 9, 10, and 11 are unconfigured. Unconfigured queues have a default scheduling weight of 1 so that they can receive a small amount of bandwidth in case they need to forward traffic. (However, queue 11 can use more of the default multdestination scheduler bandwidth if necessary to handle CPU-generated multdestination traffic.)

NOTE: Except on QFX10000 switches, all four multdestination queues, or two for QFX5200 and QFX5210, switches, have a scheduling weight of 1. Because by default multdestination traffic goes to queue 8, queue 8 receives almost all of the multdestination bandwidth. (There is no default traffic on queue 9 and queue 10, and very little default traffic on queue 11, so there is almost no competition for multdestination bandwidth.)

However, if you explicitly configure queue 9, 10, or 11 (by mapping code points to the unconfigured multdestination forwarding classes using the multdestination classifier), the explicitly configured queues share the multdestination scheduler bandwidth equally with default queue 8, because all of the queues have the same scheduling weight (1). To ensure that multdestination bandwidth is allocated to each queue properly and that the bandwidth allocation to the default queue (8) is not reduced too much, we strongly recommend that you configure a scheduler if you explicitly classify traffic into queue 9, 10, or 11.

If you map traffic to an unconfigured queue, the queue receives only the amount of group bandwidth proportional to its default weight (1). The actual amount of bandwidth an unconfigured queue receives depends on how much bandwidth the other queues in the group are using.

On QFX 10000 switches, if you map traffic to an unconfigured queue and do not schedule port resources for the queue (configure a scheduler, map it to the forwarding class that is mapped to the queue, and apply the scheduler mapping to the port), the queue receives only the amount of excess bandwidth proportional to its default weight (1). The actual amount of bandwidth an unconfigured queue gets depends on how much bandwidth the other queues on the port are using.

If the other queues use less than their allocated amount of bandwidth, the unconfigured queues can share the unused bandwidth. Configured queues have higher priority for bandwidth than unconfigured queues, so if a configured queue needs more bandwidth, then less bandwidth is available for unconfigured queues. Unconfigured queues always receive a minimum amount of bandwidth based on their scheduling weight (1). If you map traffic to an unconfigured queue, to allocate bandwidth to that queue, configure a scheduler for the forwarding class that is mapped to the queue and apply it to the port.

Default Scheduler Maps

Table 13 on page 46 shows the default mapping of forwarding classes to schedulers.

Table 13: Default Scheduler Maps

Forwarding Class	Scheduler
best-effort	Default BE scheduler
fcoe	Default FCoE scheduler
no-loss	No-loss scheduler
network-control	Default network-control scheduler
(Excluding QFX10000) mcast-be	Default multidestination scheduler

Default Shared Buffer Configuration

Table 14 on page 46 and Table 15 on page 47 show the default shared buffer allocations:

NOTE: Shared buffers do not apply to QFX10000 switches.

Table 14: Default Ingress Shared Buffer Configuration

Total Shared Ingress Buffer	Lossless Buffer	Lossless-Headroom Buffer	Lossy Buffer
100%	9%	45%	46%

Table 15: Default Egress Shared Buffer Configuration

Total Shared Egress Buffer	Lossless Buffer	Lossy Buffer	Multicast Buffer
100%	50%	31%	19%

RELATED DOCUMENTATION

<i>Overview of Junos OS CoS</i>
<i>Understanding Junos CoS Components</i>
<i>Understanding Default CoS Scheduling and Classification</i>
<i>Understanding CoS Classifiers</i>
Understanding CoS Classifiers
<i>Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces</i>
<i>Understanding CoS Code-Point Aliases</i>
<i>Understanding CoS Forwarding Classes</i>
<i>Understanding CoS Rewrite Rules</i>
<i>Understanding CoS Output Queue Schedulers</i>
<i>Understanding CoS Port Schedulers on QFX Switches</i>
<i>Understanding CoS WRED Drop Profiles</i>

CoS Inputs and Outputs Overview

Some CoS components map one set of values to another set of values. Each mapping contains one or more inputs and one or more outputs. When you configure a mapping, you set the outputs for a given set of inputs, as shown in [Table 16 on page 48](#).

Table 16: CoS Mappings—Inputs and Outputs

CoS Mappings	Inputs	Outputs	Comments
<code>classifiers</code>	<code>code-points</code>	<code>forwarding-class</code> , <code>loss-priority</code>	The map sets the forwarding class and packet loss priority (PLP) for a specific set of code points.
<code>drop-profile-map</code>	<code>loss-priority</code> , <code>protocol</code>	<code>drop-profile</code>	The map sets the drop profile for a specific PLP and protocol type.
<code>rewrite-rules</code>	<code>loss-priority</code> , <code>forwarding-class</code>	<code>code-points</code>	The map sets the code points for a specific forwarding class and PLP.
<code>rewrite-value</code> (Fibre Channel Interfaces)	<i><code>forwarding-class</code></i>	<i><code>code-point</code></i>	(Systems that support native Fibre Channel interfaces only) The map sets the code point for the forwarding class specified in the fixed classifier attached to the native Fibre Channel (NP_Port) interface.

RELATED DOCUMENTATION

| *Understanding CoS Packet Flow*

CoS Upgrade Requirements and Feature Change Compatibility

IN THIS CHAPTER

- Overview of CoS Upgrade Requirements (Junos OS Release 11.1 or 11.2 to a Later Release) | 49
- Overview of CoS Upgrade Requirements to Junos OS Release 12.2 | 51
- Overview of CoS Upgrade Requirements to Junos OS Release 12.3 (QFX3500 and QFX3600 Switches) or to Junos OS Release 13.1 (QFabric Systems) | 53
- Overview of CoS Changes Introduced in Junos OS Release 11.3 | 57
- Overview of CoS Changes Introduced in Junos OS Release 12.2 | 67

Overview of CoS Upgrade Requirements (Junos OS Release 11.1 or 11.2 to a Later Release)

Before you upgrade to Junos OS Release 11.3, you must deactivate the CoS configuration if the CoS configuration includes any of the following features:

- excess-rate option
- strict-high or high priority queues
- Any of the Junos OS Release 11.1 or 11.2 default multidestination forwarding classes



CAUTION: If your CoS configuration contains any of the features listed above and you attempt to upgrade from Junos OS Release 11.1 or 11.2 to a later version without first editing the configuration, the Junos OS might not restart.

Junos OS Release 11.3 and later for QFX Series no longer supports the excess-rate statement, the strict priority option, or the default multidestination forwarding classes used in Junos OS Release 11.1 and 11.2. In addition, Junos OS Release 11.3 introduces new restrictions on how to configure and use strict-high priority queues.

This topic does not describe how to perform the software upgrade procedure. It describes how to deactivate your CoS configuration, edit your CoS configuration, and reactivate your CoS configuration at the appropriate times.

Use the following procedure to upgrade safely from Junos OS Release 11.1 or 11.2 to a later release:

1. Deactivate the CoS configuration *before* you upgrade the software:

```
user@switch# deactivate class-of-service
```

2. Follow the upgrade procedure to Junos OS Release 11.3 or later software.
3. Make the following changes to the CoS configuration while the CoS configuration is still deactivated:
 - Remove the `excess-rate` statement from the CoS configuration if you have used it at the `[edit class-of-service schedulers]` or `[edit class-of-service traffic-control-profiles]` hierarchy level.
 - Remove the `strict-high` and `strict priority` queue configurations if you have used them at the `[edit class-of-service schedulers]` hierarchy level.
 - Remove the default multdestination forwarding classes (`mcast-be`, `mcast-af`, `mcast-ef`, and `mcast-nc`) if you have used them at the `[edit class-of-service schedulers]`, `[edit class-of-service rewrite-rules]`, `[edit class-of-service classifiers]`, `[edit class-of-service scheduler-maps]`, or `[edit class-of-service forwarding-class-sets]` hierarchy level. Alternatively, you can change the mapping of the multdestination traffic to use the new default multdestination forwarding class (`mcast`).
4. If desired, configure `strict-high` priority queues in accordance with the Junos OS Release 11.3 or later configuration rules, and map multdestination traffic to the default multdestination forwarding class (`mcast`).
5. Activate the CoS configuration:

```
user@switch# activate class-of-service
```

6. Commit the CoS configuration:

```
user@switch# commit
```

NOTE: If you configured the `transmit-rate` option for any queues under the `[edit class-of-service schedulers]` hierarchy level, if the rate is configured as an exact rate in Mbps, we recommend that

you reconfigure the `transmit-rate` option as a percentage. This is because the scheduler converts exact rates to percentages, and when the exact rate is below 1 Gbps, some granularity may be lost in the conversion. You can avoid this potential issue by specifying the `transmit-rate` option as a percentage.

RELATED DOCUMENTATION

Installing Software Packages on QFX Series Devices

Understanding CoS Classifiers

Understanding CoS Output Queue Schedulers

Understanding CoS Traffic Control Profiles

[Overview of CoS Upgrade Requirements to Junos OS Release 12.2 | 51](#)

[Overview of CoS Upgrade Requirements to Junos OS Release 12.3 \(QFX3500 and QFX3600 Switches\) or to Junos OS Release 13.1 \(QFabric Systems\) | 53](#)

Example: Configuring Unicast Classifiers

Example: Configuring Queue Schedulers

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Overview of CoS Upgrade Requirements to Junos OS Release 12.2

Before you upgrade to Junos OS Release 12.2, you might need to edit the class-of-service (CoS) configuration, because the way the QFX Series handles lossless forwarding classes has changed in Junos OS Release 12.2.

By default, the `fcoe` and `no-loss` forwarding classes are mapped to output queue 3 and output queue 4, respectively. These are the only two forwarding classes (and the only two queues) that support lossless transport.

In Junos OS Release 12.1 and earlier, explicitly setting the lossless `fcoe` and `no-loss` forwarding classes resulted in the same CoS behavior as using the default configuration. However, in Junos OS Release 12.2, the behavior when you explicitly configure the lossless forwarding classes differs from the behavior when you use the default forwarding classes.

NOTE: The default behavior differs from the explicit configuration behavior even if the explicit configuration is exactly the same as the default configuration.

- If you use the default forwarding class configuration for the lossless queues (the configuration does not include explicit setting of the `fcoe` or the `no-loss` forwarding classes), then the `fcoe` and `no-loss` queues behave as lossless queues.

If your CoS configuration does not explicitly configure the `fcoe` and `no-loss` forwarding classes, you can upgrade from Junos OS Release 12.1 to Junos OS Release 12.2, and the behavior of the two lossless queues remains the lossless.

- If your configuration includes statements that explicitly configure the `fcoe` or the `no-loss` forwarding class (using the `[set class-of-service forwarding-classes class class-name queue-num queue-number]` statement), after you upgrade to Junos OS Release 12.2, those queues do *not* receive lossless treatment and behave as lossy (best-effort) queues.

If your CoS configuration explicitly configures the `fcoe` and `no-loss` forwarding classes, to retain the lossless behavior of those queues, you need to remove the explicit configuration for these two forwarding classes from the CoS configuration *before* you upgrade.

If you upgrade to Junos OS Release 12.2 and the `fcoe` and `no-loss` forwarding classes are explicitly configured, then those two queues continue to be used, but the traffic is treated as lossy traffic, not lossless traffic. To make the queues for these two forwarding classes lossless, you must delete the explicit forwarding class configuration.



CAUTION: If you explicitly configured the `fcoe` or the `no-loss` forwarding class and you upgrade to Junos OS Release 12.2, the system does not return an upgrade error or a commit error, or a generate a syslog message, to notify you that these forwarding classes are no longer lossless. Traffic mapped to these forwarding classes is not treated as lossless traffic until you remove the explicit forwarding class configuration.

Before you upgrade, delete the `fcoe` and `no-loss` forwarding classes from the explicit configuration to preserve the lossless behavior of traffic mapped to these forwarding classes.

- To delete the explicit `fcoe` forwarding class configuration:

```
[edit]
user@switch# delete class-of-service forwarding-class class fcoe queue-num 3
user@switch# commit
```

- To delete the explicit no-loss forwarding class configuration:

```
[edit]
user@switch# delete class-of-service forwarding-class class no-loss queue-num 4
user@switch# commit
```

NOTE: If you try to delete these forwarding classes and they have not been explicitly configured on the system, the system returns the message `warning: statement not found`. This simply means that there is no explicit configuration to delete and does not change the lossless behavior of the fcoe and no-loss forwarding classes.

After you delete the explicit configuration for the fcoe and no-loss forwarding classes, traffic mapped to those forwarding classes retains its lossless behavior after the upgrade to Junos OS Release 12.2.

RELATED DOCUMENTATION

[Overview of CoS Upgrade Requirements \(Junos OS Release 11.1 or 11.2 to a Later Release\) | 49](#)

[Overview of CoS Upgrade Requirements to Junos OS Release 12.3 \(QFX3500 and QFX3600 Switches\) or to Junos OS Release 13.1 \(QFabric Systems\) | 53](#)

[Overview of CoS Changes Introduced in Junos OS Release 11.3 | 57](#)

[Installing Software Packages on QFX Series Devices](#)

[Understanding CoS Forwarding Classes](#)

[Example: Configuring Forwarding Classes](#)

Overview of CoS Upgrade Requirements to Junos OS Release 12.3 (QFX3500 and QFX3600 Switches) or to Junos OS Release 13.1 (QFabric Systems)

IN THIS SECTION

- [Support for Six Lossless Forwarding Classes | 54](#)
- [Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 56](#)

Before you upgrade to Junos OS Release 12.3 (QFX3500 and QFX3600 switches) or to Junos OS Release 13.1 (QFabric systems), you might need to edit the class-of-service (CoS) configuration, because the way the QFX Series handles lossless forwarding classes has changed from earlier Junos OS releases. (Throughout this document, changes introduced on standalone switches in Junos OS Release 12.3 are introduced on QFabric systems in Junos OS Release 13.1 unless otherwise noted.)

Support for Six Lossless Forwarding Classes

By default, the *fcoe* and *no-loss* forwarding classes are mapped to output queue 3 and output queue 4, respectively, and to IEEE 802.1p priority 3 (code point 011) and priority 4 (code point 100), respectively. These are the only two forwarding classes (and the only two queues) that support lossless transport in the default configuration.

If you use the default CoS configuration, you do not need to edit the CoS configuration after upgrading to Junos OS Release 12.3 (QFX3500 and QFX3600 switches) or to Junos OS Release 13.1 (QFabric system) because the default CoS configuration is backward-compatible.

Junos OS Release 12.3 increases the support for lossless forwarding classes (priorities) from two forwarding classes to six forwarding classes. To support configuring lossless forwarding classes, Junos OS Release 12.3 introduces a new option to forwarding class configuration: the *no-loss* packet drop attribute.

NOTE: The new *no-loss* packet drop attribute and the previously existing *no-loss* default forwarding class have the same name, but they are not the same. You can use the no-loss packet drop attribute on any unicast forwarding class.

If you explicitly configure any lossless forwarding class (including explicitly configuring the default *fcoe* and *no-loss* forwarding classes), you *must* specify the no-loss packet drop attribute to obtain lossless behavior. If you do not explicitly configure the *fcoe* and *no-loss* forwarding classes, those forwarding classes remain lossless.

The addition of the no-loss packet drop attribute to forwarding class configuration means that when you upgrade from an earlier release to Junos OS Release 12.3, the new software might not preserve the lossless forwarding class configuration of the *fcoe* and *no-loss* forwarding classes.

If you used the default forwarding class configuration for the *fcoe* and *no-loss* forwarding classes, the CoS configuration is backward-compatible. You do not have to do anything to preserve the lossless

behavior of traffic that uses those forwarding classes when you upgrade to Junos OS Release 12.3. (This is because the default configuration of these two forwarding classes includes the no-loss packet drop attribute.)

However, if you explicitly configured the fcoe or the no-loss forwarding class by including the `set forwarding-classes class forwarding-class-name queue-num queue-number` at the `[edit class-of-service]` hierarchy level, then those forwarding classes are no longer lossless, they are lossy. In Junos OS Release 12.3 and later, you must include the *no-loss* packet drop attribute in any explicit forwarding class configuration to configure a lossless forwarding class.

For example, before Junos OS Release 12.3, the following explicit configuration resulted in a lossless forwarding class:

```
user@switch# set class-of-service forwarding-classes class fcoe queue-num 3
```

However, in Junos OS Release 12.3, this configuration is lossy because it does not include the no-loss packet drop attribute. To preserve lossless behavior, after upgrading to Junos OS Release 12.3, you need to add the no-loss drop attribute:

```
user@switch# set class-of-service forwarding-classes class fcoe queue-num 3 no-loss
```

Alternatively, you can delete the explicit configuration before you upgrade to Junos OS Release 12.3 so that the system uses the default forwarding class, which is lossless:

```
user@switch# delete class-of-service forwarding-classes class fcoe queue-num 3
```

NOTE: The explicit configuration of other forwarding classes does not affect the lossless (or lossy) state of the fcoe and no-loss forwarding classes, because only the fcoe and no-loss forwarding classes are lossless forwarding classes before Junos OS Release 12.3. For example, if you explicitly configured the best-effort forwarding class but you used the default fcoe and no-loss forwarding classes in Junos OS Release 12.2, then when you upgrade to Junos OS Release 12.3, the fcoe and no-loss forwarding classes are still lossless (and the best-effort forwarding class retains its explicit configuration).

NOTE: To achieve lossless behavior for the traffic belonging to any forwarding class, you must also enable PFC on the IEEE 802.1p priority mapped to the forwarding class and ensure that DCBX exchanges the protocol TLVs for the application with the connected peer.

Scheduling on QFabric System Node Device Fabric (fte) Ports

Junos OS Release 13.1 introduces the ability to configure scheduling on the fabric (fte) ports of QFabric system Node devices. In earlier Junos OS releases, Node device fabric port scheduling was done by default, with no user configuration.

In Junos OS Release 13.1, the default fabric port scheduler configuration is similar to the default scheduler configuration on access interfaces. Similar to the access port default configuration, the default fabric port scheduler supports the five default forwarding classes (best-effort, fcoe, no-loss, network-control, and mcast). If you configure any new forwarding classes, you must configure scheduling on the fabric ports to allocate bandwidth to those forwarding classes, just as you must configure scheduling on the access ports for user-defined forwarding classes.

Strict-High Priority Scheduling on QFabric System Node Device Fabric (fte) Ports

If a fabric interface handles strict-high priority traffic, you must define a separate fc-set (priority group) for strict-high priority traffic. Strict-high priority traffic cannot be mixed with traffic of other priorities in an fc-set. For example, you might choose to create different fc-sets for best effort, lossless, strict-high priority, and multidestination traffic.

RELATED DOCUMENTATION

[Overview of CoS Upgrade Requirements \(Junos OS Release 11.1 or 11.2 to a Later Release\) | 49](#)

[Overview of CoS Upgrade Requirements to Junos OS Release 12.2 | 51](#)

Installing Software Packages on QFX Series Devices

Overview of CoS Changes Introduced in Junos OS Release 11.3

IN THIS SECTION

- [CoS Default Value Changes | 57](#)
- [Queue Priority Configuration Changes | 64](#)
- [Minimum Guaranteed Bandwidth \(Transmit Rate and Guaranteed Rate\) Changes | 65](#)
- [Excess Rate Statement Disabled | 65](#)
- [Queue Scheduling \(Low and Strict-High Priority Queues\) | 66](#)
- [Multidestination Traffic Changes | 66](#)

Junos OS Release 11.3 introduces many changes to class-of-service (CoS) functionality and to the CoS default values. This overview summarizes the changes, which other documents describe in detail.

NOTE: Some of the CoS changes are not backward compatible with Junos OS Releases 11.1 and 11.2. "[Overview of CoS Upgrade Requirements \(Junos OS Release 11.1 or 11.2 to a Later Release\)](#)" on [page 49](#) describes how to upgrade to Junos OS Release 11.3 if you have configured CoS on your QFX3500 switch.

This topic describes the following changes in CoS default values and behavior:

CoS Default Value Changes

The default values of the following CoS components have changed in Junos OS Release 11.3:

Default Forwarding Classes

In Junos OS Releases 11.1 and 11.2, there were eight default forwarding classes, four unicast default forwarding classes and four default multidestination (multicast, broadcast, and destination lookup fail) forwarding classes. [Table 17 on page 58](#) shows the old default forwarding classes and default queue mapping:

Table 17: Junos OS Release 11.1 and 11.2 Default Forwarding Classes and Queue Mapping

Default Forwarding Class	Description	Default Queue Mapping
best-effort (be)	Unicast best-effort traffic	0
no-loss	Unicast guaranteed delivery for TCP no-loss traffic	2
fcoe	Unicast guaranteed delivery for FCoE traffic	3
network-control	Unicast network control traffic	7
multicast-best-effort (mcast-be)	Multidestination best-effort traffic	8
multicast-expedited-forwarding (mcast-ef)	Multidestination low-loss, low-latency traffic	9
multicast-assured-forwarding (mcast-af)	Multidestination assured forwarding traffic	10
multicast-network-control (mcast-nc)	Multidestination network control traffic	11

Junos OS Release 11.3 changes the default forwarding classes and queue mapping in the following ways:

- Instead of eight default forwarding classes, there are five default forwarding classes.
- The same four unicast default forwarding classes remain valid, but the default queue mapping of the no-loss forwarding class has changed from queue 2 to queue 4.
- There is now only one default multidestination forwarding class instead of four default multidestination forwarding classes. All multidestination traffic is assigned by default to the default multidestination forwarding class.

NOTE: The rest of the forwarding class characteristics remain the same as before. For example, the QFX Series still supports 12 forwarding classes and 12 output queues. You can still configure a total of eight unicast forwarding classes and four multideestination forwarding classes. The unicast queues are still queues 0 through 7 and the multideestination queues are still queues 8 through 11. Unicast traffic must be mapped to unicast queues, and multideestination traffic must be mapped to multideestination queues. The queue to which a forwarding class is mapped determines whether the forwarding class is unicast or multideestination.

[Table 18 on page 59](#) shows the default forwarding classes and queue mapping in Junos OS 11.3 and later:

Table 18: Junos OS Release 11.3 Default Forwarding Classes and Queue Mapping

Default Forwarding Class	Description	Default Queue Mapping
best-effort (be)	Best-effort traffic class	0
fcoe	Guaranteed delivery for FCoE traffic	3
no-loss	Guaranteed delivery for TCP no-loss traffic	4
network-control (nc)	Network control traffic	7
mcast	Multicast traffic	8

Default IEEE 802.1p Unicast Classifiers

In Junos OS Release 11.1 and 11.2, there were default unicast classifiers only for best-effort and network-control traffic, as shown in [Table 19 on page 60](#):

Table 19: Junos OS Release 11.1 and 11.2 Default IEEE 802.1 Unicast Classifiers

Code Point	Forwarding Class	Loss Priority
be (000)	best-effort	low
be1 (001)	best-effort	low
ef (010)	best-effort	low
ef1 (011)	best-effort	low
af11 (100)	best-effort	low
af12 (101)	best-effort	low
nc1 (110)	network-control	low
nc2 (111)	network-control	low

Junos OS Release 11.3 introduces new default classifiers for FCoE and no-loss traffic, replacing the best-effort classifiers mapped to IEEE 802.1p code points 011 and 100, respectively, as shown in [Table 20 on page 60](#):

Table 20: Junos OS Release 11.3 Default IEEE 802.1 Unicast Classifiers

Code Point	Forwarding Class	Loss Priority
be (000)	best-effort	low
be1 (001)	best-effort	low
ef (010)	best-effort	low
ef1 (011)	fcoe	low

Table 20: Junos OS Release 11.3 Default IEEE 802.1 Unicast Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
af11 (100)	no-loss	low
af12 (101)	best-effort	low
nc1 (110)	network-control	low
nc2 (111)	network-control	low

Default IEEE 802.1p Multidestination Classifiers

In Junos OS Release 11.1 and 11.2, there were default multidestination classifiers for best-effort and network-control traffic, as shown in [Table 21 on page 61](#):

Table 21: Junos OS Release 11.1 and 11.2 Default IEEE 802.1 Multidestination Classifiers

Code Point	Forwarding Class	Loss Priority
be (000)	mcast-be	low
be1 (001)	mcast-be	low
ef (010)	mcast-be	low
ef1 (011)	mcast-be	low
af11 (100)	mcast-be	low
af12 (101)	mcast-be	low
nc1 (110)	mcast-nc	low

Table 21: Junos OS Release 11.1 and 11.2 Default IEEE 802.1 Multidestination Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
nc2 (111)	mcast-nc	low

Junos OS Release 11.3 replaces the best-effort and network-control multidestination classifiers and maps all IEEE 802.1p code points to the new default multidestination forwarding class, as shown in [Table 22 on page 62](#):

Table 22: Junos OS Release 11.3 Default IEEE 802.1 Multidestination Classifiers

Code Point	Forwarding Class	Loss Priority
be (000)	mcast	low
be1 (001)	mcast	low
ef (010)	mcast	low
ef1 (011)	mcast	low
af11 (100)	mcast	low
af12 (101)	mcast	low
nc1 (110)	mcast	low
nc2 (111)	mcast	low

Default Scheduler

In Junos OS Release 11.1 and 11.2, there were four default schedulers:

- Unicast best effort
- Unicast network control

- Multidestination best effort
- Multidestination network control

[Table 23 on page 63](#) shows the default scheduler configuration in Junos OS Release 11.1 and 11.2:

Table 23: Junos OS Release 11.1 and 11.2 Default Schedulers

Default Scheduler and Queue Number	Guaranteed Rate (Minimum Bandwidth)	Shaping Rate (Maximum Bandwidth)	Excess Rate (Extra Bandwidth Sharing)	Priority
Best-effort scheduler (queue 0)	75%	None	25%	Low
Network-control scheduler (queue 7)	5%	None	25%	Low
Best-effort multidestination scheduler (queue 8)	15%	None	25%	Low
Network-control multidestination scheduler (queue 11)	5%	None	25%	Low

Junos OS Release 11.3 replaces the four old classifiers with five new classifiers:

- Unicast best effort
- FCoE
- No loss
- Unicast network control
- Multidestination

There are now four different default unicast classifiers to provide default CoS for lossless queues (FCoE and no-loss traffic). Because there is only one default multidestination forwarding class in Junos OS Release 11.3, there is only one default multidestination classifier for all multidestination traffic. Also, the excess rate default value is removed from the scheduler because the `excess-rate` statement is no longer supported, as described elsewhere in this document. [Table 24 on page 64](#) shows the default scheduler configuration in Junos OS Releases 11.3:

Table 24: Default Schedulers

Default Scheduler and Queue Number	Guaranteed Rate (Minimum Bandwidth)	Shaping Rate (Maximum Bandwidth)	Excess Bandwidth Sharing	Priority
Best-effort scheduler (queue 0)	5%	None	5%	Low
FCoE scheduler (queue 3)	35%	None	35%	Low
No-loss scheduler (queue 4)	35%	None	35%	Low
Network-control scheduler (queue 7)	5%	None	5%	Low
Multidestination scheduler (queue 8)	20%	None	20%	Low

NOTE: The minimum guaranteed bandwidth rate also determines the amount of excess (extra) bandwidth that the queue can share. Extra bandwidth is allocated to queues in proportion to the minimum guaranteed bandwidth rate of each queue.

Queue Priority Configuration Changes

In Junos OS Release 11.1 and 11.2, you could configure strict-high priority queues with a guaranteed minimum bandwidth and configure forwarding class sets (priority groups) with a mix of low priority and strict-high priority queues. In Junos OS Release 11.3 and later, these configurations are invalid, and several other changes have also been implemented:

- Priority configuration in Junos OS Release 11.1 and 11.2 provided three priority levels: **strict-high**, **high**, and **low**. In Junos OS Release 11.3, the **high** priority option has been removed. Only the **strict-high** and **low** priority options are valid in Release 11.3.
- Minimum guaranteed bandwidth (transmit rate) is not allowed on strict-high priority queues. Minimum guaranteed bandwidth (guaranteed rate) is not allowed on forwarding class sets that contain strict-high priority queues.

- You cannot configure a multidestination queue as a strict-high priority queue. You cannot configure a queue as a strict-high priority queue if it belongs to the multidestination forwarding class set.
- Only one forwarding class set can contain strict-high priority queues. If you want to configure a strict-high priority queue, you must also configure a separate forwarding class set for the strict-high priority queue. A forwarding class set cannot contain a mixture of low priority and strict-high priority queues.

The rest of the queue priority characteristics remain the same as before. For example, you can configure only one queue as a strict-high priority queue.

NOTE: If you have configured **strict-high** or **high** priority queues in Junos OS Release 11.1 or 11.2, the changes in Release 11.3 are not backward compatible. Please read ["Overview of CoS Upgrade Requirements \(Junos OS Release 11.1 or 11.2 to a Later Release\)"](#) on page 49 before you upgrade to Release 11.3.

Minimum Guaranteed Bandwidth (Transmit Rate and Guaranteed Rate) Changes

The following restrictions have been placed on minimum guaranteed bandwidth configuration in Junos OS Release 11.3:

- You cannot configure a guaranteed minimum bandwidth (transmit rate) for strict-high priority queues.
- Queues (forwarding classes) with a configured transmit rate cannot be included in a forwarding class set that has strict-high priority queues.
- You cannot configure a guaranteed minimum bandwidth (guaranteed rate) for forwarding class sets that include strict-high priority queues.
- For transmit rates below 1 Gbps, we recommend that you configure the transmit rate as a percentage instead of as a fixed rate. This is because the system converts fixed rates into percentages and may round small fixed rates to a lower percentage. For example, a fixed rate of 350 Mbps is rounded down to 3 percent instead of 3.5 percent.

Excess Rate Statement Disabled

The excess-rate statement has been disabled in Junos OS Release 11.3. Excess rate was used to specify the way extra bandwidth was shared among queues.

The excess-rate statement was used at the [edit class-of-service schedulers] hierarchy level for queue scheduling configuration and at the [edit class-of-service traffic-control-profiles] hierarchy level for forwarding class set scheduling configuration.

In Junos OS Release 11.3, extra bandwidth sharing among queues is proportional to the minimum guaranteed bandwidth (transmit rate) of the queue. Extra bandwidth sharing among forwarding class sets (priority groups) is proportional to the minimum guaranteed bandwidth (guaranteed rate) of the forwarding class set.

NOTE: If you have configured the **excess-rate** option in Junos OS Release 11.1 or 11.2, the changes in Release 11.3 are not backward compatible. Please read "[Overview of CoS Upgrade Requirements \(Junos OS Release 11.1 or 11.2 to a Later Release\)](#)" on page 49 before you upgrade to Release 11.3.

Queue Scheduling (Low and Strict-High Priority Queues)

In Junos OS Release 11.1 and 11.2, if you configured a guaranteed minimum bandwidth (transmit rate) for low-priority queues, the low-priority queues received their guaranteed minimum bandwidth from the same bandwidth pool as the strict-high priority queue, using round-robin scheduling. Until the minimum bandwidth requirements of all queues were met, the strict-high priority queue and low-priority queues that had a guaranteed minimum bandwidth were treated equally. After the minimum bandwidth requirements of all queues were met, the strict-high priority queue received as much of the leftover bandwidth as it needed. This meant that the only way to ensure that a strict-high priority queue received all of the bandwidth it needed was not to configure a guaranteed minimum bandwidth for other queues.

In Junos OS Release 11.3 and later, queue scheduling has changed so that queues receive bandwidth in the following sequence:

1. The strict-high priority queue receives all of the bandwidth it needs before any other queue is served. The strict-high priority queue can take the full port bandwidth if necessary and can starve other queues on the port.
2. The guaranteed minimum bandwidth (transmit rate) of low-priority queues is served until the minimum is met or the queues are empty.
3. All other low-priority queues and needs that exceed the minimum bandwidth are served.

Multidestination Traffic Changes

The changes to the default forwarding classes and classifiers affects multidestination traffic handling in Junos OS Release 11.3:

- The number of default multidestination forwarding classes has been reduced from four default multidestination forwarding classes in Junos OS Release 11.1 and 11.2 to one default multidestination in Release 11.3 (see [Table 18 on page 59](#)).

- The default classifier configuration for multidestination traffic has changed so that there is now one default classifier for all multidestination traffic (see [Table 22 on page 62](#)).
- By default, all IEEE 802.1p code points map to the default multidestination forwarding class.
- The default scheduler for multidestination traffic has changed so that there is now one default scheduler for all multidestination traffic (see [Table 24 on page 64](#)).
- You cannot configure multidestination queues as strict-high priority queues and you cannot include strict-high priority queues in a forwarding class set that contains multidestination queues.

RELATED DOCUMENTATION

[Overview of CoS Upgrade Requirements \(Junos OS Release 11.1 or 11.2 to a Later Release\) | 49](#)

[Overview of CoS Changes Introduced in Junos OS Release 12.2 | 67](#)

[Understanding Default CoS Settings](#)

[Understanding Default CoS Scheduling and Classification](#)

[Understanding CoS Output Queue Schedulers](#)

Overview of CoS Changes Introduced in Junos OS Release 12.2

IN THIS SECTION

- [Lossless Forwarding Classes \(fcoe and no-loss\) | 68](#)
- [Default MTU for Headroom Buffer Calculation for Lossless Forwarding Classes | 68](#)
- [CoS for Layer 3 Physical Interfaces | 69](#)
- [DSCP IPv6 Classifiers and Rewrite Rules | 69](#)

Junos OS Release 12.2 introduces some changes to class-of-service (CoS) functionality and to the CoS default values. This overview summarizes the changes, which other documents describe in detail.

This topic describes the following changes in CoS default values and behavior:

Lossless Forwarding Classes (fcoe and no-loss)

The way the QFX Series handles lossless forwarding classes (the **fcoe** and **no-loss** forwarding classes) changes in Junos OS Release 12.2. In Junos OS Release 12.2 and in earlier releases, by default, the **fcoe** and **no-loss** forwarding classes are mapped to output queue 3 and output queue 4, respectively. These are the only two forwarding classes (and the only two queues) that support lossless transport.

In earlier releases, explicitly setting the lossless **fcoe** and **no-loss** forwarding classes resulted in the same CoS behavior as using the default configuration. However, in Junos OS Release 12.2, the behavior when you explicitly configure the lossless forwarding classes differs from the behavior when you use the default forwarding classes.

NOTE: The default behavior differs from the explicit configuration behavior even if the explicit configuration is exactly the same as the default configuration.

If you use the default forwarding class settings for the lossless queues (the configuration does not include explicit setting of the **fcoe** or the **no-loss** forwarding classes), then the **fcoe** and **no-loss** queues behave as lossless queues. When you upgrade to Junos OS Release 12.2, traffic assigned to the **fcoe** and **no-loss** queues continues to be treated as lossless traffic.

If your configuration explicitly sets the **fcoe** or the **no-loss** forwarding class (**set class-of-service forwarding-classes class *class-name* queue-num *queue-number***), after you upgrade to Junos OS Release 12.2, those queues do *not* receive lossless treatment and behave as lossy (**best-effort**) queues. To retain lossless treatment of the **fcoe** and **no-loss** queues, delete the explicit lossless forwarding class configuration before you upgrade to Junos OS Release 12.2.



CAUTION: If you explicitly configured the **fcoe** or the **no-loss** forwarding class, and you upgrade to Junos OS Release 12.2, the system does not return an upgrade error or a commit error, or a generate a syslog message, to notify you that these forwarding classes are no longer lossless. Traffic mapped to these forwarding classes is not treated as lossless traffic until you remove the explicit forwarding class configuration.

Default MTU for Headroom Buffer Calculation for Lossless Forwarding Classes

The default maximum transmission unit (MTU) the system uses for buffer headroom calculation is 2500 bytes for traffic classified into the **fcoe** forwarding class or the **no-loss** forwarding class.

In Junos OS Release 12.2, the default MTU used for buffer headroom calculation for the **fcoe** and **no-loss** forwarding classes remains 2500 bytes. However, if the buffer is filled, in Junos OS Release 12.2 you might experience commit failures.

CoS for Layer 3 Physical Interfaces

Before Junos OS Release 12.2, the QFX Series supported only Layer 2 CoS. Junos OS Release 12.2 introduces CoS support for Layer 3 traffic at the physical interface level.

If a physical Layer 3 interface has at least one *logical interface* configured on it, you can configure Layer 3 CoS for the physical interface. The CoS configured on the physical interface applies to all of the logical Layer 3 interfaces on that physical interface. The system does not support Layer 3 CoS configuration on individual Layer 3 logical interfaces.

DSCP IPv6 Classifiers and Rewrite Rules

Junos OS Release 12.2 introduces support for DSCP IPv6 classifiers and *rewrite rules*. The existing DSCP IP default classifier is now also the DSCP IPv6 default classifier.

You can configure and apply DSCP IPv6 classifiers and DSCP IPv6 rewrite rules to Layer 2 logical interfaces and to Layer 3 physical interfaces.

NOTE: DSCP IPv6 classifiers are not supported for multidestination (multicast, broadcast, and destination lookup fail) traffic.

RELATED DOCUMENTATION

[Overview of CoS Upgrade Requirements to Junos OS Release 12.2 | 51](#)

[Overview of CoS Changes Introduced in Junos OS Release 11.3 | 57](#)

Understanding Default CoS Settings

Understanding CoS Classifiers

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

2

PART

Classifying Traffic (Classifiers, Forwarding Classes, and Rewrite Rules)

Using Classifiers, Forwarding Classes, and Rewrite Rules | 71

Using Classifiers, Forwarding Classes, and Rewrite Rules

IN THIS CHAPTER

- Understanding CoS Classifiers | 72
- Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p) | 81
- Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p) | 84
- Example: Configuring Unicast Classifiers | 86
- Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers | 91
- Understanding Host Inbound Traffic Classification | 94
- Understanding Default CoS Scheduling and Classification | 95
- Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 106
- Understanding CoS Code-Point Aliases | 120
- Defining CoS Code-Point Aliases | 123
- Understanding CoS Forwarding Classes | 124
- Defining CoS Forwarding Classes | 131
- Example: Configuring Forwarding Classes | 133
- Understanding CoS Forwarding Class Sets (Priority Groups) | 140
- Defining CoS Forwarding Class Sets | 142
- Example: Configuring Forwarding Class Sets | 143
- Understanding Host Routing Engine Outbound Traffic Queues and Defaults | 148
- Changing the Host Outbound Traffic Default Queue Mapping | 151
- Understanding CoS Rewrite Rules | 152
- Defining CoS Rewrite Rules | 154
- Configuring CoS Fixed Classifier Rewrite Values for Native FC Interfaces (NP_Ports) | 156
- Troubleshooting an Unexpected Rewrite Value | 159

Understanding CoS Classifiers

IN THIS SECTION

- Interfaces and Output Queues | 73
- Output Queues for Unicast and Multidestination Traffic | 74
- Classifier Support by Type | 74
- Behavior Aggregate Classifiers | 75
- Fixed Classifiers on Ethernet Interfaces | 79
- Fixed Classifiers on Native Fibre Channel Interfaces (NP_Ports) | 80
- Multifield Classifiers | 80
- MPLS EXP Classifiers | 80
- Packet Classification for IRB Interfaces and RVIs | 81

Packet classification maps incoming packets to a particular class-of-service (CoS) servicing level. Classifiers map packets to a forwarding class and a loss priority, and they assign packets to output queues based on the forwarding class. There are three general types of classifiers:

- Behavior aggregate (BA) classifiers—DSCP and DSCP IPv6 classify IP and IPv6 traffic, EXP classifies MPLS traffic, and IEEE 802.1p classifies all other traffic. (Although this topic covers EXP classifiers, for more details, see [Understanding CoS MPLS EXP Classifiers and Rewrite Rules](#). EXP classifiers are applied only on family mpls interfaces.)
- Fixed classifiers—Fixed classifiers classify all ingress traffic on a physical interface into one forwarding class, regardless of the CoS bits in the packet header.
- Multifield (MF) classifiers—MF classifiers classify traffic based on more than one field in the packet header and take precedence over BA and fixed classifiers.

Classifiers assign incoming unicast and multidestination (multicast, broadcast, and destination lookup fail) traffic to forwarding classes, so that different classes of traffic can receive different treatment. Classification is based on CoS bits, DSCP bits, EXP bits, a forwarding class (fixed classifier), or packet headers (multifield classifiers). Each classifier assigns all incoming traffic that matches the classifier configuration to a particular forwarding class. Except on QFX10000 switches, classifiers and forwarding classes handle either unicast or multidestination traffic. You cannot mix unicast and multidestination traffic in the same classifier or forwarding class. On QFX10000 switches, a classifier can assign both unicast and multidestination traffic to the same forwarding class.

Interfaces and Output Queues

You can apply classifiers to Layer 2 *logical interface* unit 0 (but not to other logical interfaces), and to Layer 3 physical interfaces if the Layer 3 physical interface has at least one defined logical interface. Classifiers applied to Layer 3 physical interfaces are used on all logical interfaces on that physical interface. [Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces](#) describes the interaction between classifiers and interfaces in greater detail.

NOTE: On QFX10000 switches you can apply different classifiers to different Layer 3 logical interfaces. You cannot apply classifiers to physical interfaces.

You can configure both a BA classifier and an MF classifier on an interface. If you do this, the BA classification is performed first, and then the MF classification is performed. If the two classification results conflict, the MF classification result overrides the BA classification result.

You cannot configure a fixed classifier and a BA classifier on the same interface.

Except on QFX10000 switches, you can configure both a DSCP or DSCP IPv6 classifier and an IEEE 802.1p classifier on the same interface. IP traffic uses the DSCP or DSCP IPv6 classifier. All other traffic uses the IEEE classifier (except when you configure a global EXP classifier; in that case, MPLS traffic uses the EXP classifier providing that the interface is configured as `family mpls`). You can configure only one DSCP classifier on a physical interface (either one DSCP classifier or one DSCP IPv6 classifier, but not both).

On QFX10000 switches, you can configure either a DSCP or a DSCP IPv6 classifier and also an IEEE 802.1p classifier on the same interface. IP traffic uses the DSCP or DSCP IPv6 classifier. If you configure an interface as `family mpls`, then the interface uses the default MPLS EXP classifier. If you configure an MPLS EXP classifier, then all MPLS traffic on the switch uses the global EXP classifier. All other traffic uses the IEEE classifier. You can configure up to 64 EXP classifiers with up to 8 entries per classifier (one entry for each forwarding class) and apply them to logical interfaces.

Except on QFX10000 switches, although you can configure as many EXP classifiers as you want, the switch uses only one MPLS EXP classifier as a global classifier on all interfaces.

After you configure an MPLS EXP classifier, you can configure it as the global EXP classifier by including the EXP classifier at the `[edit class-of-service system-defaults classifiers exp]` hierarchy level. All switch interfaces that are configured as `family mpls` use the EXP classifier, on QFX10000 switches either the default or the global EXP classifier, specified in this configuration statement to classify MPLS traffic.

Output Queues for Unicast and Multidestination Traffic

NOTE: This section applies to switches except QFX10000.

You can create unicast BA classifiers for unicast traffic and multicast BA classifiers for multidestination traffic, which includes multicast, broadcast, and destination lookup fail (DLF) traffic. You cannot assign unicast traffic and multidestination traffic to the same BA classifier.

On each interface, the switch has separate output queues for unicast traffic and for multidestination traffic:

NOTE: QFX5200 switches support 10 output queues, with 8 queues dedicated to unicast traffic and 2 queues dedicated to multidestination traffic.

- The switch supports 12 output queues, with 8 queues dedicated to unicast traffic and 4 queues dedicated to multidestination traffic.
- Queues 0 through 7 are unicast traffic queues. You can apply only unicast BA classifiers to unicast queues. A unicast BA classifier should contain only forwarding classes that are mapped to unicast queues.
- Queues 8 through 11 are multidestination traffic queues. You can apply only multidestination BA classifiers to multidestination queues. A multidestination BA classifier should contain only forwarding classes that are mapped to multidestination queues.

You can apply unicast classifiers to one or more interfaces. Multidestination classifiers and EXP classifiers apply to all of the switch interfaces and cannot be applied to individual interfaces. Use the DSCP multidestination classifier for both IP and IPv6 multidestination traffic. The DSCP IPv6 classifier is not supported for multidestination traffic.

Classifier Support by Type

NOTE: This section applies only to QFX10000 switches.

You can configure enough classifiers to handle most, if not all, network scenarios. [Table 25 on page 75](#) shows how many of each type of classifiers you can configure, and how many entries you can configure per classifier.

Table 25: Classifier Support by Classifier Type

Classifier Type	Default Classifier Name	Maximum Number of Classifiers	Maximum Number of Entries per Classifier
IEEE 802.1p (Layer 2)	ieee8021p-default (for ports in trunk mode) ieee8021p-untrust (for ports in access mode)	64	16
DSCP (Layer 3)	dscp-default	64	64
DSCP IPv6 (Layer 3)	dscp-ipv6-default	64	64
EXP (MPLS)	exp-default	64	8
Fixed	There is no default fixed classifier	8	16

The number of fixed classifiers supported (8) equals the number of supported forwarding classes (fixed classifiers assign all incoming traffic on an interface to one forwarding class).

Behavior Aggregate Classifiers

Behavior aggregate classifiers map a class-of-service (CoS) value to a forwarding class and loss priority. The forwarding class determines the output queue. A scheduler uses the loss priority to control packet discard during periods of congestion by associating different drop profiles with different loss priorities.

The switch supports three types of BA classifiers:

- Differentiated Services code point (DSCP) for IP DiffServ (IP and IPv6)
- IEEE 802.1p CoS bits
- MPLS EXP (applies only to interfaces configured as `family mpls`)

BA classifiers are based on fixed-length fields, which makes them computationally more efficient than MF classifiers. Therefore, core devices, which handle high traffic volumes, are normally configured to perform BA classification.

Unicast and multicast traffic cannot share the same classifier. You can map unicast traffic and multicast traffic to the same classifier CoS value, but the unicast traffic must belong to a unicast classifier and the multicast traffic must belong to a multidestination classifier.

Default Behavior Aggregate Classification

Juniper Networks Junos OS automatically assigns implicit default classifiers to all logical interfaces based on the type of interface. [Table 26 on page 76](#) lists different types of interfaces and the corresponding implicit default BA classifiers.

Table 26: Default BA Classification

Type of Interface	Default BA Classification
Layer 2 interface in trunk mode or, except on QFX10000, tagged-access mode	ieee8021p-default
(QFX10000 only) Layer 2 interface in access mode	ieee8021p-untrusted
Layer 3 interface	dscp-default dscp-ipv6-default
(Except QFX10000) Layer 2 interface in access mode	ieee8021p-untrusted
(QFX10000 only) MPLS interface	exp-default

NOTE: Default BA classifiers assign traffic only to the best-effort, fcoe, no-loss, network-control, and, except on QFX10000 switches, mcast forwarding classes.

NOTE: Except on QFX10000 switches, there is no default MPLS EXP classifier. You must configure an EXP classifier and apply it globally to all interfaces that are configured as `family mpls` by including it in the `[edit class-of-service system-defaults classifiers exp]` hierarchy. On `family mpls` interfaces, if a fixed classifier is present on the interface, the EXP classifier overrides the fixed classifier.

If an EXP classifier is not configured, then if a fixed classifier is applied to the interface, the MPLS traffic uses the fixed classifier. If no EXP classifier and no fixed classifier is applied to the

interface, MPLS traffic is treated as best-effort traffic. DSCP classifiers are not applied to MPLS traffic.

Because the EXP classifier is global, you cannot configure some ports to use a fixed IEEE 802.1p classifier for MPLS traffic on some interfaces and the global EXP classifier for MPLS traffic on other interfaces. When you configure a global EXP classifier, all MPLS traffic on all interfaces uses the EXP classifier, even interfaces that have a fixed classifier.

When you explicitly associate a classifier with a logical interface, you override the default classifier with the explicit classifier. For other than QFX10000 switches, this applies to unicast classifiers.

NOTE: You can apply only one DSCP and one IEEE 802.1p classifier to a Layer 2 interface. If both types of classifiers are present, DSCP classifiers take precedence over IEEE 802.1p classifiers. If on QFX10000 switches you configure an EXP classifier, or on other switches a global EXP classifier, and apply it on interfaces configured as `family mpls`, then MPLS traffic uses that classifier on those interfaces.

Importing a Classifier

You can use any existing classifier, including the default classifiers, as the basis for defining a new classifier. You accomplish this using the `import` statement.

The imported classifier is used as a template and is not modified. The modifications you make become part of a new classifier (and a new template) identified by the name of the new classifier. Whenever you commit a configuration that assigns a new forwarding class-name and loss-priority value to a code-point alias or set of bits, it replaces the old entry in the new classifier template. As a result, you must explicitly specify every CoS value in every packet classification that requires modification.

Multidestination Classifiers

NOTE: This section applies to switches except QFX10000.

Multidestination classifiers are applied to all interfaces and cannot be applied to individual interfaces. You can configure both a DSCP multidestination classifier and an IEEE multidestination classifier. IP and IPv6 traffic use the DSCP classifier, and all other traffic uses the IEEE classifier.

DSCP IPv6 multidestination classifiers are not supported, so IPv6 traffic uses the DSCP multidestination classifier.

The default multdestination classifier is the IEEE 802.1p multdestination classifier.

PFC Priorities

The eight IEEE 802.1p code points correspond to the eight priorities that *priority-based flow control* (PFC) uses to differentiate traffic classes for lossless transport. When you map a forwarding class (which maps to an output queue) to an IEEE 802.1p CoS value, the IEEE 802.1p CoS value identifies the PFC priority.

Although you can map a priority to any output queue (by mapping the IEEE 802.1p code point value to a forwarding class), we recommend that the priority and the forwarding class (unicast except for QFX10000 switches) match in a one-to-one correspondence. For example, priority 0 is assigned to queue 0, priority 1 is assigned to queue 1, and so on, as shown in [Table 27 on page 78](#). A one-to-one correspondence of queue and priority numbers makes it easier to configure and maintain the mapping of forwarding classes to priorities and queues.

Table 27: Default IEEE 802.1p Code Point to PFC Priority, Output Queue, and Forwarding Class Mapping

IEEE 802.1p Code Point	PFC Priority	Output Queue (Unicast except for QFX10000)	Forwarding Class and Packet Drop Attribute
000	0	0	best-effort (drop)
001	1	1	best-effort (drop)
010	2	2	best-effort (drop)
011	3	3	fcoe (no-loss)
100	4	4	no-loss (no-loss)
101	5	5	best-effort (drop)
110	6	6	network-control (drop)
111	7	7	network-control (drop)

NOTE: By convention, deployments with converged server access typically use IEEE 802.1p priority 3 (011) for FCoE traffic. The default mapping of the `fcoe` forwarding class is to queue 3. Apply priority-based flow control (PFC) to the entire FCoE data path to configure the end-to-end lossless behavior that FCoE requires. We recommend that you use priority 3 for FCoE traffic unless your network architecture requires that you use a different priority.

Fixed Classifiers on Ethernet Interfaces

Fixed classifiers map all traffic on a physical interface to a forwarding class and a loss priority, unlike BA classifiers, which map traffic into multiple different forwarding classes based on the IEEE 802.1p CoS bits field value in the VLAN header or the DSCP field value in the type-of-service bits in the packet IP header. Each forwarding class maps to an output queue. However, when you use a fixed classifier, regardless of the CoS or DSCP bits, all Incoming traffic is classified into the forwarding class specified in the fixed classifier. A scheduler uses the loss priority to control packet discard during periods of congestion by associating different drop profiles with different loss priorities.

You cannot configure a fixed classifier and a DSCP or IEEE 802.1p BA classifier on the same interface. If you configure a fixed classifier on an interface, you cannot configure a DSCP or an IEEE classifier on that interface. If you configure a DSCP classifier, an IEEE classifier, or both classifiers on an interface, you cannot configure a fixed classifier on that interface.

NOTE: For MPLS traffic on the same interface, you can configure both a fixed classifier and an EXP classifier on QFX10000, or a global EXP classifier on other switches. When both an EXP classifier or global EXP classifier and a fixed classifier are applied to an interface, MPLS traffic on interfaces configured as `family mpls` uses the EXP classifier, and all other traffic uses the fixed classifier.

To switch from a fixed classifier to a BA classifier, or to switch from a BA classifier to a fixed classifier, deactivate the existing classifier attachment on the interface, and then attach the new classifier to the interface.

NOTE: If you configure a fixed classifier that classifies all incoming traffic into the `fcoe` forwarding class (or any forwarding class designed to handle FCoE traffic), you must ensure that all traffic that enters the interface is FCoE traffic and is tagged with the FCoE IEEE 802.1p code point (priority).

Fixed Classifiers on Native Fibre Channel Interfaces (NP_Ports)

NOTE: This section applies to switches except QFX10000.

Applying a fixed classifier to a native Fibre Channel (FC) interface (NP_Port) is a special case. By default, native FC interfaces classify incoming traffic from the FC SAN into the `fcoe` forwarding class and map the traffic to IEEE 802.1p priority 3 (code point 011). When you apply a fixed classifier to an FC interface, you also configure a priority rewrite value for the interface. The FC interface uses the priority rewrite value as the IEEE 802.1p tag value for all incoming packets instead of the default value of 3.

For example, if you specify a priority rewrite value of 5 (code point 101) for an FC interface, the interface tags all incoming traffic from the FC SAN with priority 5 and classifies the traffic into the forwarding class specified in the fixed classifier.

NOTE: The forwarding class specified in the fixed classifier on FC interfaces must be a lossless forwarding class.

Multifield Classifiers

Multifield classifiers examine multiple fields in a packet such as source and destination addresses and source and destination port numbers of the packet. With MF classifiers, you set the forwarding class and loss priority of a packet based on *firewall filter* rules.

MF classification is normally performed at the network edge because of the general lack of DiffServ code point (DSCP) support in end-user applications. On a switch at the edge of a network, an MF classifier provides the filtering functionality that scans through a variety of packet fields to determine the forwarding class for a packet. Typically, a classifier performs matching operations on the selected fields against a configured value.

MPLS EXP Classifiers

You can configure up to 64 EXP classifiers for MPLS traffic and apply them to `family mpls` interfaces. On QFX10000 switches you can use the default MPLS EXP, but on other switches there is no default MPLS classifier. You can configure an EXP classifier and apply it globally to all interfaces that are configured as `family mpls` by including it in the `[edit class-of-service system-defaults classifiers exp]` hierarchy level. On `family mpls` interfaces, if a fixed classifier is present on the interface, the EXP classifier overrides the fixed classifier for MPLS traffic only.

Except on QFX10000 switches, if an EXP classifier is not configured, then if a fixed classifier is applied to the interface, the MPLS traffic uses the fixed classifier. If no EXP classifier and no fixed classifier is

applied to the interface, MPLS traffic is treated as best-effort traffic. DSCP classifiers are not applied to MPLS traffic.

Because the EXP classifier is global, you cannot configure some ports to use a fixed IEEE 802.1p classifier for MPLS traffic on some interfaces and the global EXP classifier for MPLS traffic on other interfaces. When you configure a global EXP classifier, all MPLS traffic on all interfaces uses the EXP classifier, even interfaces that have a fixed classifier.

For details about EXP classifiers, see [Understanding CoS MPLS EXP Classifiers and Rewrite Rules](#). EXP classifiers are applied only on family `mpls` interfaces.

Packet Classification for IRB Interfaces and RVIs

On QFX10000 switches, you cannot apply classifiers directly to integrated routing and bridging (*IRB*) interfaces. Similarly, on other switches you cannot apply classifiers directly to routed VLAN interfaces (*RVIs*). This results because the members of IRBs and RVIs are VLANs, not ports. However, you can apply classifiers to the VLAN port members of an IRB interface. You can also apply MF classifiers to IRBs and RVIs.

RELATED DOCUMENTATION

Understanding CoS MPLS EXP Classifiers and Rewrite Rules
Understanding CoS Packet Flow
Understanding Default CoS Settings
Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces
Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)
Example: Configuring Unicast Classifiers
Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)
Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers
Configuring a Global MPLS EXP Classifier
Configuring Rewrite Rules for MPLS EXP Classifiers

Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)

Overview

Packet classification associates incoming packets with a particular CoS servicing level. Behavior aggregate (BA) classifiers examine the Differentiated Services code point (DSCP or DSCP IPv6) value,

the IEEE 802.1p CoS value, or the MPLS EXP value in the packet header to determine the CoS settings applied to the packet. (See [Configuring a Global MPLS EXP Classifier](#) to learn how to define EXP classifiers for MPLS traffic.) BA classifiers allow you to set the forwarding class and loss priority of a packet based on the incoming CoS value.

On most devices, unicast traffic uses different classifiers than multdestination (multicast, broadcast, and destination lookup fail) traffic. You use the `multi-destination` statement at the `[edit class-of-service]` hierarchy level to configure a multdestination BA classifier.

Multdestination classifiers apply to all of the switch interfaces and handle multicast, broadcast, and destination lookup fail (DLF) traffic. You cannot apply a multdestination classifier to a single interface or to a range of interfaces.

Platform-specific Information

- OCX Series switches do not support MPLS EXP classifiers.
- On QFX10000 switches and NFX Series devices, unicast and multdestination traffic use the same classifiers and forwarding classes.
- QFX5130, QFX5700 & QFX5220 switches do not support DSCP IPv6 classifiers and rewrite rules. However, you can apply DSCP classifiers and rewrite rules for IPV6 traffic as well.

Configuring BA Classifiers

To configure a DSCP, DSCP IPv6, or IEEE 802.1p BA classifier using the CLI:

1. Create a BA classifier:

- To create a DSCP, DSCP IPv6, or IEEE 802.1p BA classifier based on the default classifier, import the default DSCP, DSCP IPv6, or IEEE 802.1p classifier and associate it with a forwarding class, a loss priority, and a code point:

```
[edit class-of-service classifiers]
user@switch# set (dscp | dscp-ipv6 | ieee-802.1) classifier-name import default forwarding-
class forwarding-class-name loss-priority level code-points [aliases] [bit-patterns]
```

- To create a BA classifier that is not based on the default classifier, create a DSCP, DSCP IPv6, or IEEE 802.1p classifier and associate it with a forwarding class, a loss priority, and a code point:

```
[edit class-of-service classifiers]
user@switch# set (dscp | dscp-ipv6 | ieee-802.1) classifier-name forwarding-class
forwarding-class-name loss-priority level code-points [aliases] [bit-patterns]
```

2. For multidestination traffic, except on QFX10000 switches or NFX Series devices, configure the classifier as a multidestination classifier:

```
[edit class-of-service]
user@switch# set multi-destination classifiers (dscp | dscp-ipv6 | ieee-802.1 | inet-
precedence) classifier-name
```

3. Apply the classifier to a specific Ethernet interface or to all Ethernet interfaces, or to all Fibre Channel interfaces on the device.
 - To apply the classifier to a specific interface:

```
[edit class-of-service interfaces]
user@switch# set interface-name unit unit classifiers (dscp | dscp-ipv6 | ieee-802.1)
classifier-name
```

- To apply the classifier to all Ethernet interfaces on the switch, use wildcards for the interface name and the logical interface (unit) number:

```
[edit class-of-service interfaces]
user@switch# set xe-* unit * classifiers (dscp | dscp-ipv6 | ieee-802.1) classifier-name
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Unicast Classifiers

Configuring a Global MPLS EXP Classifier

Configuring Rewrite Rules for MPLS EXP Classifiers

Monitoring CoS Classifiers

Understanding CoS Classifiers

Understanding CoS Classifiers

Understanding CoS MPLS EXP Classifiers and Rewrite Rules

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)

Overview

Packet classification associates incoming packets with a particular CoS servicing level. Behavior aggregate (BA) classifiers examine the Differentiated Services code point (DSCP or DSCP IPv6) value, the IEEE 802.1p CoS value, or the MPLS EXP value in the packet header to determine the CoS settings applied to the packet. (See [Configuring a Global MPLS EXP Classifier](#) to learn how to define EXP classifiers for MPLS traffic.) BA classifiers allow you to set the forwarding class and loss priority of a packet based on the incoming CoS value.

On most devices, unicast traffic uses different classifiers than multidestination (multicast, broadcast, and destination lookup fail) traffic. You use the multi-destination statement at the [edit class-of-service] hierarchy level to configure a multidestination BA classifier.

Multidestination classifiers apply to all of the switch interfaces and handle multicast, broadcast, and destination lookup fail (DLF) traffic. You cannot apply a multidestination classifier to a single interface or to a range of interfaces.

Platform-specific Information

- OCX Series switches do not support MPLS EXP classifiers.
- On QFX10000 switches and NFX Series devices, unicast and multidestination traffic use the same classifiers and forwarding classes.
- QFX5130, QFX5700 & QFX5220 switches do not support DSCP IPv6 classifiers and rewrite rules. However, you can apply DSCP classifiers and rewrite rules for IPV6 traffic as well.

Configuring BA Classifiers

To configure a DSCP, DSCP IPv6, or IEEE 802.1p BA classifier using the CLI:

1. Create a BA classifier:

- To create a DSCP, DSCP IPv6, or IEEE 802.1p BA classifier based on the default classifier, import the default DSCP, DSCP IPv6, or IEEE 802.1p classifier and associate it with a forwarding class, a loss priority, and a code point:

```
[edit class-of-service classifiers]
user@switch# set (dscp | dscp-ipv6 | ieee-802.1) classifier-name import default forwarding-
class forwarding-class-name loss-priority level code-points [aliases] [bit-patterns]
```

- To create a BA classifier that is not based on the default classifier, create a DSCP, DSCP IPv6, or IEEE 802.1p classifier and associate it with a forwarding class, a loss priority, and a code point:

```
[edit class-of-service classifiers]
user@switch# set (dscp | dscp-ipv6 | ieee-802.1) classifier-name forwarding-class
forwarding-class-name loss-priority level code-points [aliases] [bit-patterns]
```

2. For multidestination traffic, except on QFX10000 switches or NFX Series devices, configure the classifier as a multidestination classifier:

```
[edit class-of-service]
user@switch# set multi-destination classifiers (dscp | dscp-ipv6 | ieee-802.1 | inet-
precedence) classifier-name
```

3. Apply the classifier to a specific Ethernet interface or to all Ethernet interfaces, or to all Fibre Channel interfaces on the device.

- To apply the classifier to a specific interface:

```
[edit class-of-service interfaces]
user@switch# set interface-name unit unit classifiers (dscp | dscp-ipv6 | ieee-802.1)
classifier-name
```

- To apply the classifier to all Ethernet interfaces on the switch, use wildcards for the interface name and the logical interface (unit) number:

```
[edit class-of-service interfaces]
user@switch# set xe-* unit * classifiers (dscp | dscp-ipv6 | ieee-802.1) classifier-name
```


RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Unicast Classifiers

Configuring a Global MPLS EXP Classifier

Configuring Rewrite Rules for MPLS EXP Classifiers

Monitoring CoS Classifiers

Understanding CoS Classifiers

[Understanding CoS Classifiers](#)

Understanding CoS MPLS EXP Classifiers and Rewrite Rules

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

[Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces](#)

Example: Configuring Unicast Classifiers

IN THIS SECTION

- [Requirements | 87](#)
- [Overview | 88](#)
- [Verification | 89](#)

Packet classification associates incoming packets with a particular CoS servicing level. Classifiers associate packets with a forwarding class and loss priority and assign packets to output queues based on the associated forwarding class. You apply classifiers to ingress interfaces.

Configuring Unicast Classifiers

Step-by-Step Procedure

To configure a unicast IEEE 802.1 BA classifier named **ba-ucast-classifier** as the default IEEE 802.1 map:

1. Associate code point 000 with forwarding class be and loss priority low:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 ba-ucast-classifier import default forwarding-class be loss-
priority low code-points 000
```

2. Associate code point 011 with forwarding class fcoe and loss priority low:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 ba-ucast-classifier forwarding-class fcoe loss-priority low code-
points 011
```

3. Associate code point 100 with forwarding class no-loss and loss priority low:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 ba-ucast-classifier forwarding-class no-loss loss-priority low
code-points 100
```

4. Associate code point 110 with forwarding class nc and loss priority low:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 ba-ucast-classifier forwarding-class nc loss-priority low code-
points 110
```

5. Apply the unicast classifier to ingress interface xe-0/0/10:

```
[edit class-of-service interfaces]
user@switch# set xe-0/0/10 unit 0 classifiers ieee-802.1 ba-ucast-classifier
```

Requirements

This example uses the following hardware and software components:

- One switch except QFX10000 (this example was tested on a Juniper Networks QFX3500 switch)
- Junos OS Release 11.1 or later for the QFX Series

Overview

Junos OS supports two general types of classifiers:

- Behavior aggregate or CoS value traffic classifiers—Examine the CoS value in the packet header. The value in this single field determines the CoS settings applied to the packet. BA classifiers allow you to set the forwarding class and loss priority of a packet based on the Differentiated Services code point (DSCP) value or IEEE 802.1p value.
- Multifield traffic classifiers—Examine multiple fields in the packet, such as source and destination addresses and source and destination port numbers of the packet. With multifield classifiers, you set the forwarding class and loss priority of a packet based on firewall filter rules.

NOTE: You must assign unicast traffic and multideestination (multicast, broadcast, and destination lookup fail) traffic to different classifiers. One classifier cannot include both unicast and multideestination forwarding classes. A unicast classifier can include only forwarding classes for unicast traffic.

This example describes how to configure a BA classifier called **ba-ucast-classifier** as the default IEEE 802.1 map and apply it to ingress interface **xe-0/0/10**. The BA classifier assigns loss priorities, as shown in [Table 28 on page 88](#), to incoming packets in the four forwarding classes.

You can use the same procedure to set multifield classifiers (except that you use firewall filter rules).

Table 28: ba-ucast-classifier Loss Priority Assignments

Unicast Forwarding Class	For CoS Traffic Type	ba-ucast-classifier Assignment	Packet Drop Attribute
be	Best-effort traffic	Low loss priority code point: 000	Low loss priority code point: 000
fcoe	Guaranteed delivery for Fibre Channel over Ethernet (FCoE) traffic	Low loss priority code point: 011	no-loss
no-loss	Guaranteed delivery for TCP traffic	Low loss priority code point: 100	Low loss priority code point: 100

Table 28: ba-ucast-classifier Loss Priority Assignments *(Continued)*

Unicast Forwarding Class	For CoS Traffic Type	ba-ucast-classifier Assignment	Packet Drop Attribute
nc	Network-control traffic	Low loss priority code point: 110	drop

Verification

IN THIS SECTION

- [Verifying the Unicast Classifier Configuration | 89](#)
- [Verifying the Ingress Interface Configuration | 90](#)

To verify the unicast classifier configuration, perform these tasks:

Verifying the Unicast Classifier Configuration

Purpose

Verify that you configured the unicast classifier with the correct forwarding classes, loss priorities, and code points.

Action

List the classifier configuration using the operational mode command `show configuration class-of-service classifiers ieee-802.1 ba-ucast-classifier`:

```
user@switch> show configuration class-of-service classifiers ieee-802.1 ba-ucast-classifier
  forwarding-class be {
    loss-priority low code-points 000;
  }
  forwarding-class fcoe {
    loss-priority low code-points 011;
  }
  forwarding-class no-loss {
```

```

        loss-priority low code-points 100;
    }
    forwarding-class nc
        loss-priority low code-points 110;
    }

```

Verifying the Ingress Interface Configuration

Purpose

Verify that the unicast classifier ba-ucast-classifier is attached to ingress interface xe-0/0/10.

Action

List the ingress interface using the operational mode command `show configuration class-of-service interfaces xe-0/0/10`:

```

user@switch> show configuration class-of-service interfaces xe-0/0/10
congestion-notification-profile fcoe-cnp;
unit 0 {
    classifiers {
        ieee-802.1 ba-ucast-classifier;
    }
}

```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers

Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)

Configuring a Global MPLS EXP Classifier

Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)

Monitoring CoS Classifiers

Understanding CoS Classifiers

[Understanding CoS Classifiers](#)

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers

IN THIS SECTION

- [Requirements | 92](#)
- [Overview | 92](#)
- [Verification | 93](#)

Packet classification associates incoming packets with a particular CoS servicing level. Behavior aggregate (BA) classifiers examine the CoS value in the packet header to determine the CoS settings applied to the packet. BA classifiers allow you to set the forwarding class and loss priority of a packet based on the incoming CoS value.

Beginning with Junos OS Release 17.1, EX4300 switches support multidestination classifiers. On EX4300 switches, you can apply multidestination classifiers globally or to a specific interface. If you apply multidestination classifiers both globally and to a specific interface, the classifications on the interface take precedence.

Multidestination classifiers apply to all of the switch interfaces and handle multicast, broadcast, and destination lookup fail (DLF) traffic. You cannot apply a multidestination classifier to a single interface or to a range of interfaces, except on an EX4300 switch.

Unicast and multidestination traffic must use different classifiers.

Configuring Multidestination Classifiers

Step-by-Step Procedure

To configure a multicast IEEE 802.1 BA classifier named `ba-mcast-classifier`:

1. Associate code point `000` with forwarding class `mcast` and loss priority `low`:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 ba-mcast-classifier forwarding-class mcast loss-priority low code-points 000
```

2. Configure the classifier as a multdestination classifier:

```
[edit class-of-service]
user@switch# set multi-destination classifiers ieee-802.1 ba-mcast-classifier
```

Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 11.1 or later for the QFX Series.

Overview

Junos OS supports three general types of classifiers:

- Behavior aggregate or CoS value traffic classifiers—Examine the CoS value in the packet header. The value in this single field determines the CoS settings applied to the packet. BA classifiers allow you to set the forwarding class and loss priority of a packet based on the CoS value.
- Fixed classifiers. Fixed classifiers classify all ingress traffic on a physical interface into one forwarding class, regardless of the CoS bits in the VLAN header or the DSCP bits in the packet header.
- Multifield traffic classifiers—Examine multiple fields in the packet such as source and destination addresses and source and destination port numbers of the packet. With multifield classifiers, you set the forwarding class and loss priority of a packet based on firewall filter rules.

Multidestination classifiers apply to all of the switch interfaces and handle multicast, broadcast, and destination lookup fail (DLF) traffic. You cannot apply a multidestination classifier to a single interface or to a range of interfaces.

NOTE: You must assign unicast traffic and multicast traffic to different classifiers. One classifier cannot include both unicast and multicast forwarding classes. A multidestination classifier can include only forwarding classes for multicast traffic.

The following example describes how to configure a BA classifier called `ba-mcast-classifier`, which is applied to all of the switch interfaces. The BA classifier assigns loss priorities, as shown in [Table 29 on page 93](#), to incoming packets in the multidestination forwarding class.

You can also use firewall filters to set multifield classifiers.

Table 29: BA-mcast-classifier Loss Priority Assignments

Multicast Forwarding Class	Traffic Type	ba-mcast-classifier Assignment
mcast	Best-effort multicast traffic	Low loss priority code point: 000

Verification

IN THIS SECTION

- [Verifying the IEEE 802.1 Multidestination Classifier | 93](#)
- [Verifying the Multidestination Classifier Configuration | 94](#)

To verify the multidestination classifier configuration, perform these tasks:

Verifying the IEEE 802.1 Multidestination Classifier

Purpose

Verify that the classifier `ba-mcast-classifier` is configured as the IEEE 802.1 multidestination classifier:

Action

Verify the results of the classifier configuration using the operational mode command `show configuration class-of-service multi-destination classifiers ieee-802.1`:

```
user@switch> show configuration class-of-service multi-destination classifiers ieee-802.1
ba-mcast-classifier;
```


Verifying the Multidestination Classifier Configuration

Purpose

Verify that you configured the multidestination classifier with the correct forwarding classes, loss priorities, and code points.

Action

List the classifier configuration using the operational mode command `show configuration class-of-service classifiers ieee-802.1 ba-mcast-classifier`:

```
user@switch> show configuration class-of-service classifiers ieee-802.1 ba-mcast-classifier
  forwarding-class mcast {
    loss-priority low code-points 000;
  }
```

Release History Table

Release	Description
17.1	Beginning with Junos OS Release 17.1, EX4300 switches support multidestination classifiers.

RELATED DOCUMENTATION

<i>Example: Configuring Unicast Classifiers</i>
<i>Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)</i>
<i>Monitoring CoS Classifiers</i>
<i>Understanding CoS Classifiers</i>
Understanding CoS Classifiers
<i>Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces</i>

Understanding Host Inbound Traffic Classification

The destination address of traffic that enters the switch can be an external device such as another switch, a router, or a server, or the destination can be the host (the switch Routing Engine or CPU). When the destination is an external device, the DSCP and IEEE 802.1p code-point bits of incoming

traffic are preserved as the traffic travels through the switch to the egress port. At the egress port, the code-point bits are either preserved when the packets are sent to the next hop or they are rewritten according to the rewrite rule attached to the egress interface.

When the destination of incoming traffic is the host, DSCP bits are preserved. However, IEEE 802.1p bits are not preserved. The IEEE 802.1p bits of traffic destined for the host are set to zero (0). This does not affect system behavior because the switch prioritizes traffic destined for the host based on the protocol type. For example, the switch gives a higher priority to BPDU traffic than to ping traffic.

Understanding Default CoS Scheduling and Classification

IN THIS SECTION

- [Default Classification | 96](#)
- [Default Scheduling | 101](#)
- [Default DCBX Advertisement | 105](#)
- [Default Scheduling and Classification Summary | 105](#)

If you do not explicitly configure classifiers and apply them to interfaces, the switch uses the default classifier to group ingress traffic into forwarding classes. If you do not configure scheduling on an interface, the switch uses the default schedulers to provide egress port resources for traffic. Default classification maps all traffic into default forwarding classes (best-effort, fcoe, no-loss, network-control, and mcast). Each default forwarding class has a default scheduler, so that the traffic mapped to each default forwarding class receives port bandwidth, prioritization, and packet drop characteristics.

The switch supports direct port scheduling and enhanced transmission selection (ETS), also known as hierarchical port scheduling, except on QFX5200 and QFX5210 switches.

Hierarchical scheduling groups IEEE 802.1p priorities (IEEE 802.1p code points, which classifiers map to forwarding classes, which in turn are mapped to output queues) into priority groups (forwarding class sets). If you use only the default traffic scheduling and classification, the switch automatically creates a default priority group that contains all of the priorities (which are mapped to forwarding classes and output queues), and assigns 100 percent of the port output bandwidth to that priority group. The forwarding classes (queues) in the default forwarding class set receive bandwidth based on the default classifier settings. The default priority group is transparent. It does not appear in the configuration and is used for Data Center Bridging Capability Exchange (DCBX) protocol advertisement.

NOTE: If you explicitly configure one or more priority groups on an interface, any forwarding class that is not assigned to a priority group on that interface receives *no bandwidth*. This means that if you configure hierarchical scheduling on an interface, every forwarding class (priority) that you want to forward traffic on that interface must belong to a forwarding class set (priority group). ETS is not supported on QFX5200 or QFX5210 switches.

The following sections describe:

Default Classification

On switches except QFX10000 and NFX Series devices, the default classifiers assign unicast and multicast best-effort and network-control ingress traffic to default forwarding classes and loss priorities. The switch applies default unicast IEEE 802.1, unicast DSCP, and multideestination classifiers to each interface that does not have explicitly configured classifiers.

On QFX10000 switches and NFX Series devices, the default classifiers assign ingress traffic to default forwarding classes and loss priorities. The switch applies default IEEE 802.1, DSCP, and DSCP IPv6 classifiers to each interface that does not have explicitly configured classifiers. If you do not configure and apply EXP classifiers for MPLS traffic to logical interfaces, MPLS traffic on interfaces configured as family mpls uses the IEEE classifier.

If you explicitly configure one type of classifier but not other types of classifiers, the system uses only the configured classifier and does not use default classifiers for other types of traffic. There are two default IEEE 802.1 classifiers: a trusted classifier for ports that are in trunk mode or tagged-access mode, and an untrusted classifier for ports that are in access mode.

NOTE: The default classifiers apply to unicast traffic except on QFX10000 switches and NFX Series devices. Tagged-access mode does not apply to QFX10000 switches or NFX Series devices.

Table 30 on page 97 shows the default mapping of IEEE 802.1 code-point values to forwarding classes and loss priorities for ports in trunk mode or tagged-access mode.

Table 30: Default IEEE 802.1 Classifiers for Ports in Trunk Mode or Tagged-Access Mode (Trusted Classifier)

Code Point	Forwarding Class	Loss Priority
be (000)	best-effort	low
be1 (001)	best-effort	low
ef (010)	best-effort	low
ef1 (011)	fcoe	low
af11 (100)	no-loss	low
af12 (101)	best-effort	low
nc1 (110)	network-control	low
nc2 (111)	network-control	low

Table 31 on page 97 shows the default mapping of IEEE 802.1p code-point values to forwarding classes and loss priorities for ports in access mode (all incoming traffic is mapped to best-effort forwarding classes).

NOTE: Table 31 on page 97 applies only to unicast traffic except on QFX10000 switches and NFX Series devices.

Table 31: Default IEEE 802.1 Classifiers for Ports in Access Mode (Untrusted Classifier)

Code Point	Forwarding Class	Loss Priority
000	best-effort	low
001	best-effort	low

Table 31: Default IEEE 802.1 Classifiers for Ports in Access Mode (Untrusted Classifier) (Continued)

Code Point	Forwarding Class	Loss Priority
010	best-effort	low
011	best-effort	low
100	best-effort	low
101	best-effort	low
110	best-effort	low
111	best-effort	low

Table 32 on page 98 shows the default mapping of IEEE 802.1 code-point values to multideestination (multicast, broadcast, and destination lookup fail traffic) forwarding classes and loss priorities.

NOTE: Table 32 on page 98 does not apply to QFX10000 switches or NFX Series devices.

Table 32: Default IEEE 802.1 Multideestination Classifiers

Code Point	Forwarding Class	Loss Priority
be (000)	mcast	low
be1 (001)	mcast	low
ef (010)	mcast	low
ef1 (011)	mcast	low
af11 (100)	mcast	low

Table 32: Default IEEE 802.1 Multidestination Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
af12 (101)	mcast	low
nc1 (110)	mcast	low
nc2 (111)	mcast	low

Table 33 on page 99 shows the default mapping of DSCP code-point values to forwarding classes and loss priorities for DSCP IP and DCSP IPv6.

NOTE: Table 33 on page 99 applies only to unicast traffic except on QFX10000 switches and NFX Series devices.

Table 33: Default DSCP IP and IPv6 Classifiers

Code Point	Forwarding Class	Loss Priority
ef (101110)	best-effort	low
af11 (001010)	best-effort	low
af12 (001100)	best-effort	low
af13 (001110)	best-effort	low
af21 (010010)	best-effort	low
af22 (010100)	best-effort	low
af23 (010110)	best-effort	low
af31 (011010)	best-effort	low

Table 33: Default DSCP IP and IPv6 Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
af32 (011100)	best-effort	low
af33 (011110)	best-effort	low
af41 (100010)	best-effort	low
af42 (100100)	best-effort	low
af43 (100110)	best-effort	low
be (000000)	best-effort	low
cs1 (001000)	best-effort	low
cs2 (010000)	best-effort	low
cs3 (011000)	best-effort	low
cs4 (100000)	best-effort	low
cs5 (101000)	best-effort	low
nc1 (110000)	network-control	low
nc2 (111000)	network-control	low

NOTE: There are no default DSCP IP or IPv6 multdestination classifiers for multdestination traffic. DSCP IPv6 multdestination classifiers are not supported for multdestination traffic.

[Table 34 on page 101](#) shows the default mapping of MPLS EXP code-point values to forwarding classes and loss priorities, which apply only on QFX10000 switches and NFX Series devices.

Table 34: Default EXP Classifiers on QFX10000 Switches and NFX Series Devices

Code Point	Forwarding Class	Loss Priority
000	best-effort	low
001	best-effort	high
010	expedited-forwarding	low
011	expedited-forwarding	high
100	assured-forwarding	low
101	assured-forwarding	high
110	network-control	low
111	network-control	high

Default Scheduling

The default schedulers allocate egress bandwidth resources to egress traffic as shown in [Table 35 on page 102](#):

Table 35: Default Scheduler Configuration

Default Scheduler and Queue Number	Transmit Rate (Guaranteed Minimum Bandwidth)	Shaping Rate (Maximum Bandwidth)	Excess Bandwidth Sharing	Priority	Buffer Size
best-effort forwarding class scheduler (queue 0)	5% 15% (QFX10000, NFX Series)	None	5% 15% (QFX10000, NFX Series)	low	5% 15% (QFX10000, NFX Series)
fcoe forwarding class scheduler (queue 3)	35%	None	35%	low	35%
no-loss forwarding class scheduler (queue 4)	35%	None	35%	low	35%
network-control forwarding class scheduler (queue 7)	5% 15% (QFX10000, NFX Series)	None	5% 15% (QFX10000, NFX Series)	low	5% 15% (QFX10000, NFX Series)
(Excluding QFX10000 and NFX Series) mcast forwarding class scheduler (queue 8)	20%	None	20%	low	20%

NOTE: By default, the minimum guaranteed bandwidth (transmit rate) determines the amount of excess (extra) bandwidth that a queue can share. Extra bandwidth is allocated to queues in proportion to the transmit rate of each queue. On switches that support the `excess-rate` statement, you can override the default setting and configure the excess bandwidth percentage independently of the transmit rate on queues that are not strict-high priority queues.

By default, only the four (QFX10000 switches and NFX Series devices) or five (other switches) default schedulers shown in [Table 35 on page 102](#) have traffic mapped to them. Only the forwarding classes

and queues associated with the default schedulers receive default bandwidth, based on the default scheduler transmit rate. (You can configure schedulers and forwarding classes to allocate bandwidth to other queues or to change the bandwidth and other scheduling properties of a default queue.)

On QFX10000 switches and NFX Series devices, if a forwarding class does not transport traffic, the bandwidth allocated to that forwarding class is available to other forwarding classes. Unicast and multdestination (multicast, broadcast, and destination lookup fail) traffic use the same forwarding classes and output queues.

On switches other than QFX10000 and NFX Series devices, multdestination queue 11 receives enough bandwidth from the default multdestination scheduler to handle CPU-generated multdestination traffic.

On QFX10000 and NFX Series devices, default scheduling is port scheduling. Default hierarchical scheduling, known as enhanced transmission selection (ETS, defined in IEEE 802.1Qaz), allocates the total port bandwidth to the four default forwarding classes served by the four default schedulers, as defined by the four default schedulers. The result is the same as direct port scheduling. Configuring hierarchical port scheduling, however, enables you to group forwarding classes that carry similar types of traffic into forwarding class sets (also called priority groups), and to assign port bandwidth to each forwarding class set. The port bandwidth assigned to the forwarding class set is then assigned to the forwarding classes within the forwarding class set. This hierarchy enables you to control port bandwidth allocation with greater granularity, and enables hierarchical sharing of extra bandwidth to better utilize link bandwidth.

Except on QFX10000 switches and NFX Series devices, default hierarchical scheduling divides the total port bandwidth between two groups of traffic: unicast traffic and multdestination traffic. By default, unicast traffic consists of queue 0 (best-effort forwarding class), queue 3 (fcoe forwarding class), queue 4 (no-loss forwarding class), and queue 7 (network-control forwarding class). Unicast traffic receives and shares a total of 80 percent of the port bandwidth. By default, multdestination traffic (mcast queue 8) receives a total of 20 percent of the port bandwidth. So on a 10-Gigabit port, unicast traffic receives 8-Gbps of bandwidth and multdestination traffic receives 2-Gbps of bandwidth.

NOTE: Except on QFX5200, QFX5210, and QFX10000 switches and NFX Series devices, which do not support queue 11, multdestination queue 11 also receives a small amount of default bandwidth from the multdestination scheduler. CPU-generated multdestination traffic uses queue 11, so you might see a small number of packets egress from queue 11. In addition, in the unlikely case that firewall filter match conditions map multdestination traffic to a unicast forwarding class, that traffic uses queue 11.

Default scheduling uses weighted round-robin (WRR) scheduling. Each queue receives a portion (weight) of the total available interface bandwidth. The scheduling weight is based on the transmit rate of the default scheduler for that queue. For example, queue 7 receives a default scheduling weight of 5 percent, or 15 percent on QFX10000 and NFX Series devices, of the available bandwidth, and queue 4

receives a default scheduling weight of 35 percent of the available bandwidth. Queues are mapped to forwarding classes, so forwarding classes receive the default bandwidth for the queues to which they are mapped.

On QFX10000 switches and NFX Series devices, for example, queue 7 is mapped to the network-control forwarding class and queue 4 is mapped to the no-loss forwarding class. Each forwarding class receives the default bandwidth for the queue to which it is mapped. Unused bandwidth is shared with other default queues.

If you want non-default (unconfigured) queues to forward traffic, you should explicitly map traffic to those queues (configure the forwarding classes and queue mapping) and create schedulers to allocate bandwidth to those queues. By default, queues 1, 2, 5, and 6 are unconfigured.

Except on QFX5200, QFX5210, and QFX10000 switches and NFX Series devices, which do not support them, multidestination queues 9, 10, and 11 are unconfigured. Unconfigured queues have a default scheduling weight of 1 so that they can receive a small amount of bandwidth in case they need to forward traffic. However, queue 11 can use more of the default multidestination scheduler bandwidth if necessary to handle CPU-generated multidestination traffic.

NOTE: All four (two on QFX5200 and QFX5210 switches) multidestination queues have a scheduling weight of 1. Because by default multidestination traffic goes to queue 8, queue 8 receives almost all of the multidestination bandwidth. (There is no traffic on queue 9 and queue 10, and very little traffic on queue 11, so there is almost no competition for multidestination bandwidth.)

However, if you explicitly configure queue 9, 10, or 11 (by mapping code points to the unconfigured multidestination forwarding classes using the multidestination classifier), the explicitly configured queues share the multidestination scheduler bandwidth equally with default queue 8, because all of the queues have the same scheduling weight (1). To ensure that multidestination bandwidth is allocated to each queue properly and that the bandwidth allocation to the default queue (8) is not reduced too much, we strongly recommend that you configure a scheduler if you explicitly classify traffic into queue 9, 10, or 11.

If you map traffic to an unconfigured queue, the queue receives only the amount of excess bandwidth proportional to its default weight (1). The actual amount of bandwidth an unconfigured queue gets depends on how much bandwidth the other queues are using.

If some queues use less than their allocated amount of bandwidth, the unconfigured queues can share the unused bandwidth. Sharing unused bandwidth is one of the key advantages of hierarchical port scheduling. Configured queues have higher priority for bandwidth than unconfigured queues, so if a configured queue needs more bandwidth, then less bandwidth is available for unconfigured queues. Unconfigured queues always receive a minimum amount of bandwidth based on their scheduling weight (1). If you map traffic to an unconfigured queue, to allocate bandwidth to that queue, configure a scheduler for the forwarding class that is mapped to the queue.

Default DCBX Advertisement

When you configure hierarchical scheduling on an interface, DCBX advertises each priority group, the priorities in each priority group, and the bandwidth properties of each priority and priority group.

If you do not configure hierarchical scheduling on an interface, DCBX advertises the automatically created default priority group and its priorities. DCBX also advertises the default bandwidth allocation of the priority group, which is 100 percent of the port bandwidth.

Default Scheduling and Classification Summary

If you do not configure scheduling on an interface:

- Default classifiers classify ingress traffic.
- Default schedulers schedule egress traffic.
- DCBX advertises a single default priority group with 100 percent of the port bandwidth allocated to that priority group. All priorities (forwarding classes) are assigned to the default priority group and receive bandwidth based on their default schedulers. The default priority group is generated automatically and is not user-configurable.

RELATED DOCUMENTATION

Understanding CoS Packet Flow

Understanding CoS Hierarchical Port Scheduling (ETS)

Understanding Default CoS Settings

Understanding CoS Virtual Output Queues (VOQs) on QFX10000 Switches

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

Understanding DCB Features and Requirements

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

Example: Configuring Unicast Classifiers

Example: Configuring Queue Schedulers

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

IN THIS SECTION

- Supported Classifier and Rewrite Rule Types | 106
- Ethernet Interfaces Supported for Classifier and Rewrite Rule Configuration | 109
- Default Classifiers | 112
- Default Rewrite Rules | 113
- Classifier Precedence | 113
- Classifier Behavior and Limitations | 115
- Rewrite Rule Precedence and Behavior | 116
- Classifier and Rewrite Rule Configuration Interaction with Ethernet Interface Configuration | 117

At ingress interfaces, classifiers group incoming traffic into classes based on the IEEE 802.1p, DSCP, or MPLS EXP *class of service* (CoS) code points in the packet header. At egress interfaces, you can use *rewrite rules* to change (re-mark) the code point bits before the interface forwards the packets.

You can apply classifiers and rewrite rules to interfaces to control the level of CoS applied to each packet as it traverses the system and the network. This topic describes:

Supported Classifier and Rewrite Rule Types

Table 36 on page 106 shows the supported types of classifiers and rewrite rules supports:

Table 36: Supported Classifiers and Rewrite Rules

Classifier or Rewrite Rule Type	Description
Fixed classifier	Classifies all ingress traffic on a physical interface into one fixed forwarding class, regardless of the CoS bits in the packet header.
DSCP and DSCP IPv6 unicast classifiers	Classifies IP and IPv6 traffic into forwarding classes and assigns loss priorities to the traffic based on DSCP code point bits.

Table 36: Supported Classifiers and Rewrite Rules *(Continued)*

Classifier or Rewrite Rule Type	Description
IEEE 802.1p unicast classifier	Classifies Ethernet traffic into forwarding classes and assigns loss priorities to the traffic based on IEEE 802.1p code point bits.
MPLS EXP classifier	<p>Classifies MPLS traffic into forwarding classes and assigns loss priorities to the traffic on interfaces configured as family mpls.</p> <p>QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 switches, and QFabric systems, use one global EXP classifier on all family mpls switch interfaces.</p> <p>QFX10000 switches do not support global EXP classifiers. You can apply the same EXP classifier or different EXP classifiers to different family mpls interfaces.</p>
DSCP multidestination classifier (also used for IPv6 multidestination traffic) NOTE: This applies only to switches that use different classifiers for unicast and multidestination traffic. It does not apply to switches that use the same classifiers for unicast and multidestination traffic.	<p>Classifies IP and IPv6 multicast, broadcast, and destination lookup fail (DLF) traffic into multidestination forwarding classes.</p> <p>Multidestination classifiers are applied to all interfaces and cannot be applied to individual interfaces.</p>
IEEE 802.1p multidestination classifier NOTE: This applies only to switches that use different classifiers for unicast and multidestination traffic. It does not apply to switches that use the same classifiers for unicast and multidestination traffic.	<p>Classifies Ethernet multicast, broadcast, and destination lookup fail (DLF) traffic into multidestination forwarding classes.</p> <p>Multidestination classifiers are applied to all interfaces and cannot be applied to individual interfaces.</p>
DSCP and DSCP IPv6 rewrite rules	Re-marks the DSCP code points of IP and IPv6 packets before forwarding the packets.
IEEE 802.1p rewrite rule	Re-marks the IEEE 802.1p code points of Ethernet packets before forwarding the packets.

Table 36: Supported Classifiers and Rewrite Rules *(Continued)*

Classifier or Rewrite Rule Type	Description
MPLS EXP rewrite rule	Re-marks the EXP code points of MPLS packets before forwarding the packets on interfaces configured as family mpls.

NOTE: On switches that support native Fibre Channel (FC) interfaces, you can specify a rewrite value on native FC interfaces (NP_Ports) to set the IEEE 802.1p code point of incoming FC traffic when the NP_Port encapsulates the FC packet in Ethernet before forwarding it to the FCoE network (see ["Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway" on page 620](#)).

DSCP, IEEE 802.1p, and MPLS EXP classifiers are behavior aggregate (BA) classifiers. On QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, unlike DSCP and IEEE 802.1p classifiers, EXP classifiers are global and apply only to all interfaces that are configured as family mpls. On QFX10000 switches, you apply EXP classifiers to individual logical interfaces, and different interfaces can use different EXP classifiers.

Unlike DSCP and IEEE 802.1p BA classifiers, there is no default EXP classifier. Also unlike DSCP and IEEE 802.1p classifiers, for MPLS traffic on family mpls interfaces only, EXP classifiers overwrite fixed classifiers. (An interface that has a fixed classifier uses the EXP classifier for MPLS traffic, not the fixed classifier, and the fixed classifier is used for all other traffic.)

On switches that use different classifiers for unicast and multdestination traffic, multdestination classifiers are global and apply to all interfaces; you cannot apply a multdestination classifier to individual interfaces.

Classifying packets into forwarding classes assigns packets to the output queues mapped to those forwarding classes. The traffic classified into a forwarding class receives the CoS scheduling configured for the output queue mapped to that forwarding class.

NOTE: In addition to BA classifiers and fixed classifiers, which classify traffic based on the CoS field in the packet header, you can use firewall filters to configure multifield (MF) classifiers. MF classifiers classify traffic based on more than one field in the packet header and take precedence over BA and fixed classifiers.

Ethernet Interfaces Supported for Classifier and Rewrite Rule Configuration

To apply a classifier to incoming traffic or a rewrite rule to outgoing traffic, you need to apply the classifier or rewrite rule to one or more interfaces. When you apply a classifier or rewrite rule to an interface, the interface uses the classifier to group incoming traffic into forwarding classes and uses the rewrite rule to re-mark the CoS code point value of each packet before it leaves the system.

Not all interfaces types support all types of CoS configuration. This section describes:

Interface Types That Support Classifier and Rewrite Rule Configuration

You can apply classifiers and rewrite rules to Ethernet interfaces. For Layer 3 LAGs, configure BA or fixed classifiers on the LAG (ae) interface. The classifier configured on the LAG is valid on all of the LAG member interfaces.

On switches that support native FC interfaces, you can apply fixed classifiers to native FC interfaces (NP_Ports). You cannot apply other types of classifiers or rewrite rules to native FC interfaces. You can rewrite the value of the IEEE 802.1p code point of incoming FC traffic when the interface encapsulates it in Ethernet before forwarding it to the FCoE network as described in ["Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway" on page 620](#).

Classifier and Rewrite Rule Physical and Logical Ethernet Interface Support

The Ethernet ports can function as:

- Layer 2 physical interfaces (family ethernet-switching)
- Layer 2 logical interfaces (family ethernet-switching)
- Layer 3 physical interfaces (family inet/inet6)
- Layer 3 logical interfaces (family inet/inet6)
- MPLS interfaces (family mpls)

You can apply CoS classifiers and rewrite rules only to the following interfaces:

- Layer 2 logical interface

NOTE: On a Layer 2 interface, use **unit *** to apply the rule to all of the logical units on that interface.

- On QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, Layer 3 physical interfaces if at least one logical Layer 3 interface is configured on the physical interface

NOTE: The CoS you configure on a Layer 3 physical interface is applied to all of the Layer 3 logical interfaces on that physical interface. This means that each Layer 3 interface uses the same classifiers and rewrite rules for all of the Layer 3 traffic on that interface.

- On QFX10000 switches, Layer 3 logical interfaces. You can apply different classifiers and rewrite rules to different Layer 3 logical interfaces.

Ethernet Interface Support for Most QFX Series Switches, and QFabric Systems

You cannot apply classifiers or rewrite rules to Layer 2 physical interfaces or to Layer 3 logical interfaces. [Table 37 on page 110](#) shows on which interfaces you can configure and apply classifiers and rewrite rules.

NOTE: The CoS feature support listed in this table is identical on single interfaces and aggregated Ethernet interfaces.

Table 37: Ethernet Interface Support for Classifier and Rewrite Rule Configuration (QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 Switches, and QFabric Systems)

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interface (unit * applies rule to all logical interfaces)	Layer 3 Physical Interfaces (If at Least One Logical Layer 3 Interface Is Defined)	Layer 3 Logical Interfaces
Fixed classifier	No	Yes	Yes	No
DSCP classifier	No	Yes	Yes	No
DSCP IPv6 classifier	No	Yes	Yes	No
IEEE 802.1p classifier	No	Yes	Yes	No
EXP classifier	Global classifier, applies only to all switch interfaces that are configured as family mpls. Cannot be configured on individual interfaces.			

Table 37: Ethernet Interface Support for Classifier and Rewrite Rule Configuration (QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 Switches, and QFabric Systems) (Continued)

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interface (unit * applies rule to all logical interfaces)	Layer 3 Physical Interfaces (If at Least One Logical Layer 3 Interface Is Defined)	Layer 3 Logical Interfaces
DSCP rewrite rule	No	Yes	Yes	No
DSCP IPv6 rewrite rule	No	Yes	Yes	No
IEEE 802.1p rewrite rule	No	Yes	Yes	No
EXP rewrite rule	No	Yes	Yes	No

NOTE: IEEE 802.1p multdestination and DSCP multdestination classifiers are applied to all interfaces and cannot be applied to individual interfaces. No DSCP IPv6 multdestination classifier is supported. IPv6 multdestination traffic uses the DSCP multdestination classifier.

Ethernet Interface Support for QFX10000 Switches

You cannot apply classifiers or rewrite rules to Layer 2 or Layer 3 physical interfaces. You can apply classifiers and rewrite rules only to Layer 2 logical interface unit 0. You can apply different classifiers and rewrite rules to different Layer 3 logical interfaces. [Table 38 on page 112](#) shows on which interfaces you can configure and apply classifiers and rewrite rules.

NOTE: The CoS feature support listed in this table is identical on single interfaces and aggregated Ethernet interfaces.

Table 38: Ethernet Interface Support for Classifier and Rewrite Rule Configuration (QFX10000 Switches)

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interface (Unit 0 Only)	Layer 3 Physical Interfaces	Layer 3 Logical Interfaces
Fixed classifier	No	Yes	No	Yes
DSCP classifier	No	Yes	No	Yes
DSCP IPv6 classifier	No	Yes	No	Yes
IEEE 802.1p classifier	No	Yes	No	Yes
EXP classifier	No	Yes	No	Yes
DSCP rewrite rule	No	Yes	No	Yes
DSCP IPv6 rewrite rule	No	Yes	No	Yes
IEEE 802.1p rewrite rule	No	Yes	No	Yes
EXP rewrite rule	No	Yes	No	Yes

Routed VLAN Interfaces (RVIs) and Integrated Routing and Bridging (IRB) Interfaces

You cannot apply classifiers and rewrite rules directly to routed VLAN interfaces (RVIs) or integrated routing and bridging (IRB) interfaces because the members of RVIs and IRBs are VLANs, not ports. However, you can apply classifiers and rewrite rules to the VLAN port members of an *RVI* or an *IRB*. You can also apply MF classifiers to RVIs and IRBs.

Default Classifiers

If you do not explicitly configure classifiers on an Ethernet interface, the switch applies default classifiers so that the traffic receives basic CoS treatment. The factors that determine the default classifier applied

to the interface include the interface type (Layer 2 or Layer 3), the port mode (trunk, tagged-access, or access), and whether logical interfaces have been configured.

The switch applies default classifiers using the following rules:

- If the physical interface has at least one Layer 3 *logical interface* configured, the logical interfaces use the default DSCP classifier.
- If the physical interface has a Layer 2 logical interface in trunk mode or tagged-access mode, it uses the default IEEE 802.1p trusted classifier.

NOTE: Tagged-access mode is available only on QFX3500 and QFX3600 devices when used as standalone switches or as QFabric system Node devices.

- If the physical interface has a Layer 2 logical interface in access mode, it uses the default IEEE 802.1p untrusted classifier.
- If the physical interface has no logical interface configured, no default classifier is applied.
- On switches that use different classifiers for unicast and multidestination traffic, the default multidestination classifier is the IEEE 802.1p multidestination classifier.
- There is no default MPLS EXP classifier. If you want to classify MPLS traffic using EXP bits on these switches, on QFX10000 switches, configure an EXP classifier and apply it to a logical interface that is configured as `family mpls`. On QFX5100, QFX5200, EX4600, QFX3500 and QFX3600 switches, and on QFabric systems, configure an EXP classifier and configure it as the global system default EXP classifier.

Default Rewrite Rules

No default rewrite rules are applied to interfaces. If you want to re-mark packets at the egress interface, you must explicitly configure a rewrite rule.

Classifier Precedence

You can apply multiple classifiers (MF, fixed, IEEE 802.1p, DSCP, or EXP) to an Ethernet interface to handle different types of traffic. (EXP classifiers are global and apply only to all MPLS traffic on all `family mpls` interfaces.) When you apply more than one classifier to an interface, the system uses an order of precedence to determine which classifier to use on interfaces:

Classifier Precedence on Physical Ethernet Interfaces (QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 Switches, and QFabric Systems)

QFX10000 switches do not support configuring classifiers on physical interfaces. The precedence of classifiers on physical interfaces, from the highest-priority classifier to the lowest-priority classifier, is:

- MF classifier on a logical interface (no classifier has a higher priority than MF classifiers)
- Fixed classifier on the physical interface
- DSCP or DSCP IPv6 classifier on the physical interface
- IEEE 802.1p classifier on the physical interface

NOTE: If an EXP classifier is configured, MPLS traffic uses the EXP classifier on all family mpls interfaces, even if an MF or fixed classifier is applied to the interface. If an EXP classifier is not configured, then if a fixed classifier is applied to the interface, the MPLS traffic uses the fixed classifier. If no EXP classifier and no fixed classifier is applied to the interface, MPLS traffic is treated as best-effort traffic. DSCP classifiers are not applied to MPLS traffic.

You can apply a DSCP classifier, an IEEE 802.1p classifier, and an EXP classifier on a physical interface. When all three classifiers are on an interface, IP traffic uses the DSCP classifier, MPLS traffic on family mpls interfaces uses the EXP classifier, and all other traffic uses the IEEE classifier.

NOTE: You cannot apply a fixed classifier and a DSCP or IEEE classifier to the same interface. If a DSCP classifier, an IEEE classifier, or both are on an interface, you cannot apply a fixed classifier to that interface unless you first delete the DSCP and IEEE classifiers. If a fixed classifier is on an interface, you cannot apply a DSCP classifier or an IEEE classifier unless you first delete the fixed classifier.

Classifier Precedence on Logical Ethernet Interfaces (All Switches)

The precedence of classifiers on logical interfaces, from the highest priority classifier to the lowest priority classifier, is:

- MF classifier on a logical interface (no classifier has a higher priority than MF classifiers).
- Fixed classifier on the logical interface.
- DSCP or DSCP IPv6 classifier on the physical or logical interface..
- IEEE 802.1p classifier on the physical or logical interface.

NOTE: If a global EXP classifier is configured, MPLS traffic uses the EXP classifier on all `family mpls` interfaces, even if a fixed classifier is applied to the interface. If a global EXP classifier is not configured, then:

- If a fixed classifier is applied to the interface, the MPLS traffic uses the fixed classifier. If no EXP classifier and no fixed classifier is applied to the interface, MPLS traffic is treated as best-effort traffic.

You can apply both a DSCP classifier and an IEEE 802.1p classifier on a logical interface. When both a DSCP and an IEEE classifier are on an interface, IP traffic uses the DSCP classifier, and all other traffic uses the IEEE classifier. Only MPLS traffic on interfaces configured as `family mpls` uses the EXP classifier.

Classifier Behavior and Limitations

Consider the following behaviors and constraints when you apply classifiers to Ethernet interfaces. Behaviors for applying classifiers to physical interfaces do not pertain to QFX10000 switches.

- You can configure only one DSCP classifier (IP or IPv6) on a physical interface. You cannot configure both types of DSCP classifier on one physical interface. Both IP and IPv6 traffic use whichever DSCP classifier is configured on the interface.
- When you configure a DSCP or a DSCP IPv6 classifier on a physical interface and the physical interface has at least one logical Layer 3 interface, all packets (IP, IPv6, and non-IP) use that classifier.
- An interface with both a DSCP classifier (IP or IPv6) and an IEEE 802.1p classifier uses the DSCP classifier for IP and IPv6 packets, and uses the IEEE classifier for all other packets.
- Fixed classifiers and BA classifiers (DSCP and IEEE classifiers) are not permitted simultaneously on an interface. If you configure a fixed classifier on an interface, you cannot configure a DSCP or an IEEE classifier on that interface. If you configure a DSCP classifier, an IEEE classifier, or both classifiers on an interface, you cannot configure a fixed classifier on that interface.
- When you configure an IEEE 802.1p classifier on a physical interface and a DSCP classifier is not explicitly configured on that interface, the interface uses the IEEE classifier for all types of packets. No default DSCP classifier is applied to the interface. (In this case, if you want a DSCP classifier on the interface, you must explicitly configure it and apply it to the interface.)
- The system does not apply a default classifier to a physical interface until you create a logical interface on that physical interface. If you configure a Layer 3 logical interface, the system uses the default DSCP classifier. If you configure a Layer 2 logical interface, the system uses the default IEEE 802.1p trusted classifier if the port is in trunk mode or tagged-access mode, or the default IEEE 802.1p untrusted classifier if the port is in access mode.

- MF classifiers configured on logical interfaces take precedence over BA and fixed classifiers, with the exception of the global EXP classifier, which is always used for MPLS traffic on family `mpls` interfaces. (Use firewall filters to configure MF classifiers.) When BA or fixed classifiers are present on an interface, you can still configure an MF classifier on that interface.
- There is no default EXP classifier for MPLS traffic.
- You can configure up to 64 EXP classifiers. On QFX10000 switches, you can apply different EXP classifiers to different interfaces.

However, on On QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, the switch uses only one MPLS EXP classifier as a global classifier on all family `mpls` interfaces. After you configure an MPLS EXP classifier, you can configure it as the global EXP classifier by including the EXP classifier in the `[edit class-of-service system-defaults classifiers exp]` hierarchy level.

All family `mpls` switch interfaces use the EXP classifier specified using this configuration statement to classify MPLS traffic, even on interfaces that have a fixed classifier. No other traffic uses the EXP classifier.

Rewrite Rule Precedence and Behavior

The following rules apply on Ethernet interfaces for rewrite rules:

- If you configure one DSCP (or DSCP IPv6) rewrite rule and one IEEE 802.1p rewrite rule on an interface, both rewrite rules take effect. Traffic with IP and IPv6 headers use the DSCP rewrite rule, and traffic with a VLAN tag uses the IEEE rewrite rule.
- If you do not explicitly configure a rewrite rule, there is no default rewrite rule, so the system does not apply any rewrite rule to the interface.
- You can apply a DSCP rewrite rule or a DSCP IPv6 rewrite rule to an interface, but you cannot apply both a DSCP and a DSCP IPv6 rewrite rule to the same interface. Both IP and IPv6 packets use the same DSCP rewrite rule, regardless of whether the configured rewrite rule is DSCP or DSCP IPv6.
- MPLS EXP rewrite rules apply only to logical interfaces on family `mpls` interfaces. You cannot apply to an EXP rewrite rule to a physical interface. You can configure up to 64 EXP rewrite rules, but you can only use 16 EXP rewrite rules at any time on the switch.
- A logical interface can use both DSCP (or DSCP IPv6) and EXP rewrite rules.
- DSCP and DSCP IPv6 rewrite rules are not applied to MPLS traffic.
- If the switch is performing penultimate hop popping (PHP), EXP rewrite rules do not take effect. If both an EXP classifier and an EXP rewrite rule are configured on the switch, then the EXP value from

the last popped label is copied into the inner label. If either an EXP classifier or an EXP rewrite rule (but not both) is configured on the switch, then the inner label EXP value is sent unchanged.

NOTE: On each physical interface, either all forwarding classes that are being used on the interface must have rewrite rules configured or no forwarding classes that are being used on the interface can have rewrite rules configured. On any physical port, do not mix forwarding classes with rewrite rules and forwarding classes without rewrite rules.

NOTE: Rewrite rules are applied *before* the egress filter is matched to traffic. Because the code point rewrite occurs before the egress filter is matched to traffic, the egress filter match is based on the rewrite value, not on the original code point value in the packet.

Classifier and Rewrite Rule Configuration Interaction with Ethernet Interface Configuration

On QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 switches used as standalone switches or as QFabric system Node devices, you can apply classifiers and rewrite rules only on Layer 2 logical interface unit 0 and Layer 3 physical interfaces (if the Layer 3 physical interface has at least one defined logical interface). On QFX10000 switches, you can apply classifiers and rewrite rules only to Layer 2 logical interface unit 0 and to Layer 3 logical interfaces. This section focuses on BA classifiers, but the interaction between BA classifiers and interfaces described in this section also applies to fixed classifiers and rewrite rules.

NOTE: On QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 switches used as standalone switches or as QFabric system Node devices, EXP classifiers, are global and apply to all switch interfaces. See [Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\)](#) for how to configure multdestination classifiers and see [Configuring a Global MPLS EXP Classifier](#) for how to configure EXP classifiers.

On switches that use different classifiers for unicast and multdestination traffic, multdestination classifiers are global and apply to all switch interfaces.

There are two components to applying classifiers or rewrite rules to interfaces:

1. Setting the interface family (inet, inet6, or ethernet-switching; ethernet-switching is the default interface family) in the [edit interfaces] configuration hierarchy.
2. Applying a classifier or rewrite rule to the interface in the [edit class-of-service] hierarchy.

These are separate operations that can be set and committed at different times. Because the type of classifier or rewrite rule you can apply to an interface depends on the interface family configuration, the system performs checks to ensure that the configuration is valid. The method the system uses to notify you of an invalid configuration depends on the set operation that causes the invalid configuration.

NOTE: QFX10000 switches cannot be misconfigured in the following two ways because you can configure classifiers only on logical interfaces. Only switches that allow classifier configuration on physical and logical interfaces can experience the following misconfigurations.

If applying the classifier or rewrite rule to the interface in the [edit class-of-service] hierarchy causes an invalid configuration, the system rejects the configuration and returns a commit check error.

If setting the interface family in the [edit interfaces] configuration hierarchy causes an invalid configuration, the system creates a syslog error message. If you receive the error message, you need to remove the classifier or rewrite rule configuration from the logical interface and apply it to the physical interface, or remove the classifier or rewrite rule configuration from the physical interface and apply it to the logical interface. For classifiers, if you do not take action to correct the error, the system programs the default classifier for the interface family on the interface. (There are no default rewrite rules. If the commit check fails, no rewrite rule is applied to the interface.)

Two scenarios illustrate these situations:

- Applying a classifier to an Ethernet interface causes a commit check error
- Configuring the Ethernet interface family causes a syslog error

These scenarios differ on different switches because some switches support classifiers on physical Layer 3 interfaces but not on logical Layer 3 interfaces, while other switches support classifiers on logical Layer 3 interfaces but not on physical Layer 3 interfaces.

Two scenarios illustrate these situations:

NOTE: Both of these scenarios also apply to fixed classifiers and rewrite rules.

QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 Switch Scenarios

The following scenarios also apply the QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 switches when they are used as QFabric system Node devices.

Scenario 1: Applying a Classifier to an Ethernet Interface Causes a Commit Check Error

In Scenario 1, we set the interface family, and then specify an invalid classifier.

1. Set and commit the interface as a Layer 3 (family inet) interface:

```
[edit interfaces]
user@switch# set xe-0/0/20 unit 0 family inet
user@switch# commit
```

This commit operation succeeds.

2. Set and commit a DSCP classifier on the logical interface (this example uses a DSCP classifier named dscp1):

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 unit 0 classifiers dscp dscp1
user@switch# commit
```

This configuration is not valid, because it attempts to apply a classifier to a Layer 3 logical interface. Because the failure is caused by the class-of-service configuration and not by the interface configuration, the system rejects the commit operation and issues a commit error, not a syslog message.

Note that the commit operation succeeds if you apply the classifier to the physical Layer 3 interface as follows:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 classifiers dscp dscp1
user@switch# commit
```

Because the logical unit is not specified, the classifier is applied to the physical Layer 3 interface in a valid configuration, and the commit check succeeds.

Scenario 2: Configuring the Ethernet Interface Family Causes a Syslog Error

In Scenario 2, we set the classifier first, and then set an invalid interface type.

1. Set and commit a DSCP classifier on a logical interface that has no existing configuration:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 unit 0 classifiers dscp dscp1
user@switch# commit
```

This commit succeeds. Because no explicit configuration existed on the interface, it is by default a Layer 2 (family ethernet-switching) interface. Layer 2 logical interfaces support BA classifiers, so applying the classifier is a valid configuration.

2. Set and commit the interface as a Layer 3 interface (family inet) interface:

```
[edit interfaces]
user@switch# set xe-0/0/20 unit 0 family inet
user@switch# commit
```

This configuration is not valid because it attempts to change an interface from Layer 2 (family ethernet-switching) to Layer 3 (family inet) when a classifier has already been applied to a logical interface. Layer 3 logical interfaces do not support classifiers. Because the failure is caused by the interface configuration and not by the class-of-service configuration, the system does not issue a commit error, but instead issues a syslog message.

When the system issues the syslog message, it programs the default classifier for the interface type on the interface. In this scenario, the interface has been configured as a Layer 3 interface, so the system applies the default DSCP profile to the physical Layer 3 interface.

In this scenario, to install a configured DSCP classifier, remove the misconfigured classifier from the Layer 3 logical interface and apply it to the Layer 3 physical interface. For example:

```
[edit]
user@switch# delete class-of-service interfaces xe-0/0/20 unit 0 classifiers dscp dscp1
user@switch# commit
user@switch# set class-of-service interfaces xe-0/0/20 classifiers dscp dscp1
user@switch# commit
```

RELATED DOCUMENTATION

Understanding CoS Packet Flow
Configuring CoS

Understanding CoS Code-Point Aliases

A code-point alias assigns a name to a pattern of code-point bits. You can use this name instead of the bit pattern when you configure other CoS components such as classifiers and *rewrite rules*.

NOTE: This topic applies to all EX Series switches except the EX4600. Because the EX4600 uses a different chipset than other EX Series switches, the code-point aliases on EX4600 match those on QFX Series switches. For EX4600 code-point aliases, see [Understanding CoS Code-Point Aliases](#).

Behavior aggregate classifiers use class-of-service (CoS) values such as Differentiated Services Code Points (DSCPs) or IEEE 802.1 bits to associate incoming packets with a particular forwarding class and the CoS servicing level associated with that forwarding class. You can assign a meaningful name or alias to the CoS values and use that alias instead of bits when configuring CoS components. These aliases are not part of the specifications but are well known through usage. For example, the alias for DSCP 101110 is widely accepted as ef (expedited forwarding).

When you configure forwarding classes and define classifiers, you can refer to the markers by alias names. You can configure code point alias names for user-defined classifiers. If the value of an alias changes, it alters the behavior of any classifier that references it.

You can configure code-point aliases for the following type of CoS markers:

- dscp or dscp-ipv6—Handles incoming IP and IPv6 packets.
- ieee-802.1—Handles Layer 2 frames.

[Table 39 on page 121](#) shows the default mapping of code-point aliases to IEEE code points.

Table 39: Default IEEE 802.1 Code-Point Aliases

CoS Value Types	Mapping
be	000
be1	001
ef	010
ef1	011
af11	100
af12	101

Table 39: Default IEEE 802.1 Code-Point Aliases (Continued)

CoS Value Types	Mapping
nc1	110
nc2	111

[Table 40 on page 122](#) shows the default mapping of code-point aliases to DSCP and DSCP IPv6 code points.

Table 40: Default DSCP and DSCP IPv6 Code-Point Aliases

CoS Value Types	Mapping
ef	101110
af11	001010
af12	001100
af13	001110
af21	010010
af22	010100
af23	010110
af31	011010
af32	011100
af33	011110
af41	100010

Table 40: Default DSCP and DSCP IPv6 Code-Point Aliases *(Continued)*

CoS Value Types	Mapping
af42	100100
af43	100110
be	000000
cs1	001000
cs2	010000
cs3	011000
cs4	100000
cs5	101000
nc1	110000
nc2	111000

RELATED DOCUMENTATION

Understanding Junos CoS Components

Defining CoS Code-Point Aliases

Defining CoS Code-Point Aliases

You can use code-point aliases to streamline the process of configuring CoS features on your switch. A code-point alias assigns a name to a pattern of code-point bits. You can use this name instead of the bit pattern when you configure other CoS components such as classifiers and rewrite rules.

You can configure code-point aliases for the following CoS marker types:

- DSCP or DSCP IPv6—Handles incoming IPv4 or IPv6 packets.
- IEEE 802.1p—Handles Layer 2 frames.

To configure a code-point alias:

1. Specify a CoS marker type (IEEE 802.1 or DSCP).
2. Assign an alias.
3. Specify the code point that corresponds to the alias.

```
[edit class-of-service code-point-aliases]
user@switch# set (dscp | dscp-ipv6 | ieee-802.1) alias-name code-point-bits
```

For example, to configure a code-point alias for an IEEE 802.1 CoS marker type that has the alias name be2 and maps to the code-point bits 001:

```
[edit class-of-service code-point-aliases]
user@switch# set ieee-802.1 be2 001
```

RELATED DOCUMENTATION

Monitoring CoS Code-Point Value Aliases

Understanding CoS Code-Point Aliases

Understanding CoS Forwarding Classes

IN THIS SECTION

- Default Forwarding Classes | 126
- Forwarding Class Configuration Rules | 127
- Lossless Transport Support | 129

Forwarding classes group traffic and assign the traffic to output queues. Each forwarding class is mapped to an output queue. Classification maps incoming traffic to forwarding classes based on the code point bits in the packet or frame header. Forwarding class to queue mapping defines the output queue used for the traffic classified into a forwarding class.

Except on NFX Series devices, a classifier must associate each packet with one of the following four (QFX10000 switches) or five (other switches) default forwarding classes or with a user-configured forwarding class to assign an output queue to the packet:

- fcoe—Guaranteed delivery for Fibre Channel over Ethernet (FCoE) traffic.
- no-loss—Guaranteed delivery for TCP lossless traffic.
- best-effort—Provides best-effort delivery without a service profile. Loss priority is typically not carried in a class-of-service (CoS) value.
- network-control—Supports protocol control and is typically high priority.
- mcast—(Except QFX10000) Delivery of multdestination (multicast, broadcast, and destination lookup fail) packets.

On NFX Series devices, a classifier must associate each packet with one of the following four default forwarding classes or with a user-configured forwarding class to assign an output queue to the packet:

- best-effort (be)—Provides no service profile. Loss priority is typically not carried in a CoS value.
- expedited-forwarding (ef)—Provides a low loss, low latency, low jitter, assured bandwidth, end-to-end service.
- assured-forwarding (af)—Provides a group of values you can define and includes four subclasses: AF1, AF2, AF3, and AF4, each with two drop probabilities: low and high.
- network-control (nc)—Supports protocol control and thus is typically high priority.

The switch supports up to eight (QFX10000 and NFX Series devices), 10 (QFX5200 switches), or 12 (other switches) forwarding classes, thus enabling flexible, differentiated, packet classification. For example, you can configure multiple classes of best-effort traffic such as **best-effort**, **best-effort1**, and **best-effort2**.

On QFX10000 and NFX Series devices, unicast and multdestination (multicast, broadcast, and destination lookup fail) traffic use the same forwarding classes and output queues.

Except on QFX10000 and NFX Series devices, a switch supports 8 queues for unicast traffic (queues 0 through 7) and 2 (QFX5200 switches) or 4 (other switches) output queues for multdestination traffic (queues 8 through 11). Forwarding classes mapped to unicast queues are associated with unicast traffic, and forwarding classes mapped to multdestination queues are associated with multdestination traffic. You cannot map unicast and multdestination traffic to the same queue. You cannot map a strict-high

priority queue to a multdestination forwarding class because queues 8 through 11 do not support strict-high priority configuration.

Default Forwarding Classes

Table 41 on page 126 shows the four default forwarding classes that apply to all switches but not NFX Series devices. Except on QFX10000, these forwarding classes apply to unicast traffic. You can rename the forwarding classes. Assigning a new forwarding class name does not alter the default classification or scheduling applied to the queue that is mapped to that forwarding class. CoS configurations can be complex, so unless it is required by your scenario, we recommend that you use the default class names and queue number associations.

Table 41: Default Forwarding Classes

Forwarding Class Name	Default Queue Mapping	Comments
best-effort	0	<p>The software does not apply any special CoS handling to best-effort traffic. This is a backward compatibility feature. Best-effort traffic is usually the first traffic to be dropped during periods of network congestion.</p> <p>By default, this is a lossy forwarding class with a packet drop attribute of drop.</p>
fcoe	3	<p>By default, the fcoe forwarding class is a lossless forwarding class designed to handle Fibre Channel over Ethernet (FCoE) traffic. The no-loss packet drop attribute is applied by default.</p> <p>NOTE: By convention, deployments with converged server access typically use IEEE 802.1p priority 3 (011) for FCoE traffic. The default mapping of the fcoe forwarding class is to queue 3. Apply <i>priority-based flow control</i> (PFC) to the entire FCoE data path to configure the end-to-end lossless behavior that FCoE requires.</p> <p>We recommend that you use priority 3 for FCoE traffic unless your network architecture requires that you use a different priority.</p>
no-loss	4	<p>By default, this is a lossless forwarding class with a packet drop attribute of no-loss.</p>

Table 41: Default Forwarding Classes *(Continued)*

Forwarding Class Name	Default Queue Mapping	Comments
network-control	7	<p>The software delivers packets in this service class with a high priority. (These packets are not delay-sensitive.)</p> <p>Typically, these packets represent routing protocol hello or keepalive messages. Because loss of these packets jeopardizes proper network operation, packet delay is preferable to packet discard.</p> <p>By default, this is a lossy forwarding class with a packet drop attribute of drop.</p>

NOTE: [Table 42 on page 127](#) applies only to multidestination traffic except on QFX10000 switches and NFX Series devices.

Table 42: Default Forwarding Classes for Multidestination Packets

Forwarding Class Name	Default Queue Mapping	Comments
mcast	8	<p>The software does not apply any special CoS handling to the multidestination packets. These packets are usually dropped under congested network conditions.</p> <p>By default, this is a lossy forwarding class with a packet drop attribute of drop.</p>

NOTE: Mirrored traffic is always sent to the queue that corresponds to the multidestination forwarding class. The switched copy of the mirrored traffic is forwarded with the priority determined by the behavior aggregate classification process.

Forwarding Class Configuration Rules

Take the following rules into account when you configure forwarding classes:

Queue Assignment Rules

The following rules govern queue assignment:

- CoS configurations that specify more queues than the switch can support are not accepted. The commit operation fails with a detailed message that states the total number of queues available.
- All default CoS configurations are based on queue number. The name of the forwarding class that appears in the default configuration is the forwarding class currently mapped to that queue.
- (Except QFX10000 and NFX Series devices) Only unicast forwarding classes can be mapped to unicast queues (0 through 7), and only multdestination forwarding classes can be mapped to multdestination queues (8 through 11).
- (Except QFX10000 and NFX Series devices) Strict-high priority queues cannot be mapped to multdestination forwarding classes. (Strict-high priority traffic cannot be mapped to queues 8 through 11).
- If you map more than one forwarding class to a queue, all of the forwarding classes mapped to the same queue must have the same packet drop attribute: either all of the forwarding classes must be lossy or all of the forwarding classes must be lossless.

You can limit the amount of traffic that receives strict-high priority treatment on a strict-high priority queue by configuring a transmit rate. The transmit rate sets the amount of traffic on the queue that receives strict-high priority treatment. The switch treats traffic that exceeds the transmit rate as low priority traffic that receives the queue excess rate bandwidth. Limiting the amount of traffic that receives strict-high priority treatment prevents other queues from being starved while also ensuring that the amount of traffic specified in the transmit rate receives strict-high priority treatment.

NOTE: Except on QFX10000 and NFX Series devices, you can use the *shaping-rate* statement to throttle the rate of packet transmission by setting a maximum bandwidth. On QFX10000 and NFX Series devices, you can use the transmit rate to set a limit on the amount of bandwidth that receives strict-high priority treatment on a strict-high priority queue.

On QFX10000 and NFX Series devices, if you configure more than one strict-high priority queue on a port, you must configure a transmit rate on each of the strict-high priority queues. If you configure more than one strict-high priority queue on a port and you do not configure a transmit rate on the strict-high priority queues, the switch treats only the first queue you configure as a strict-high priority queue. The switch treats the other queues as low priority queues. If you configure a transmit rate on some strict-high priority queues but not on other strict-high priority queues on a port, the switch treats the queues that have a transmit rate as strict-high priority queues, and treats the queues that do not have a transmit rate as low priority queues.

Scheduling Rules

When you configure a forwarding class and map traffic to it (that is, you are not using a default classifier and forwarding class), you must also define a scheduling policy for the forwarding class.

Defining a scheduling policy means:

- Mapping a scheduler to the forwarding class in a scheduler map
- Including the forwarding class in a forwarding class set
- Associating the scheduler map with a traffic control profile
- Attaching the traffic control profile to a forwarding class set and applying the traffic control profile to an interface

On QFX10000 switches and NFX Series devices, you can define a scheduling policy using port scheduling as follows:

- Mapping a scheduler to the forwarding class in a scheduler map
- Applying the scheduler map to one or more interfaces

Rewrite Rules

On each physical interface, either all forwarding classes that are being used on the interface must have rewrite rules configured, or no forwarding classes that are being used on the interface can have rewrite rules configured. On any physical port, do not mix forwarding classes with rewrite rules and forwarding classes without rewrite rules.

Lossless Transport Support

The switch supports up to six lossless forwarding classes. For lossless transport, you must enable PFC on the IEEE 802.1p code point of lossless forwarding classes. The following limitations apply to support lossless transport:

- The external cable length from the switch or QFabric system Node device to other devices cannot exceed 300 meters.
- The internal cable length from the QFabric system Node device to the QFabric system Interconnect device cannot exceed 150 meters.
- For FCoE traffic, the interface maximum transmission unit (MTU) must be at least 2180 bytes to accommodate the packet payload, headers, and checks.
- Changing any portion of a PFC configuration on a port blocks the entire port until the change is completed. After a PFC change is completed, the port is unblocked and traffic resumes. Changing the

PFC configuration means any change to a congestion notification profile that is configured on a port (enabling or disabling PFC on a code point, changing the MRU or cable-length value, or specifying an output flow control queue). Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

NOTE: QFX10002-60C does not support PFC and lossless queues; that is, default lossless queues (fcoe and no-loss) will be lossy queues.

NOTE: Junos OS Release 12.2 introduces changes to the way lossless forwarding classes (the fcoe and no-loss forwarding classes) are handled.

In Junos OS Release 12.1, both explicitly configuring the fcoe and no-loss forwarding classes, and using the default configuration for these forwarding classes, resulted in the same lossless behavior for traffic mapped to those forwarding classes.

However, in Junos OS Release 12.2, if you explicitly configure the fcoe or the no-loss forwarding class, that forwarding class is no longer treated as a lossless forwarding class. Traffic mapped to these forwarding classes is treated as lossy (best-effort) traffic. This is true even if the explicit configuration is exactly the same as the default configuration.

If your CoS configuration from Junos OS Release 12.1 or earlier includes the explicit configuration of the fcoe or the no-loss forwarding class, then when you upgrade to Junos OS Release 12.2, those forwarding classes are not lossless. To preserve the lossless treatment of these forwarding classes, delete the explicit fcoe and no-loss forwarding class configuration before you upgrade to Junos OS Release 12.2.

See "[Overview of CoS Changes Introduced in Junos OS Release 12.2](#)" on page 67 for detailed information about this change and how to delete an existing lossless configuration.

In Junos OS Release 12.3, the default behavior of the fcoe and no-loss forwarding classes is the same as in Junos OS Release 12.2. However, in Junos OS Release 12.3, you can configure up to six lossless forwarding classes. All explicitly configured lossless forwarding classes must include the new no-loss packet drop attribute or the forwarding class is lossy.

RELATED DOCUMENTATION

[Overview of CoS Changes Introduced in Junos OS Release 12.2 | 67](#)

Understanding Junos CoS Components

Understanding CoS Packet Flow

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Example: Configuring Forwarding Classes

Defining CoS Forwarding Classes

Defining CoS Forwarding Classes

Forwarding classes allow you to group packets for transmission. The switch supports a total of eight (QFX10000 and NFX Series devices), 10 (QFX5200 switches), or 12 (other switches) forwarding classes. To forward traffic, you map (assign) the forwarding classes to output queues. Starting in Junos OS Release 22.1R1, QFX10000 Series devices support 16 forwarding classes.

The QFX10000 switches and NFX Series devices have eight output queues, queues 0 through 7. These queues support both unicast and multdestination traffic.

Except on QFX10000 and NFX Series devices, the switch has 10 output queues (QFX5200) or 12 output queues (other switches). Queues 0 through 7 are for unicast traffic and queues 8 through 11 are for multicast traffic. Forwarding classes mapped to unicast queues must carry unicast traffic, and forwarding classes mapped to multdestination queues must carry multdestination traffic. There are four default unicast forwarding classes and one default multdestination forwarding class.

The default forwarding classes, except on NFX Series devices, are:

NOTE: Except on QFX10000, these are the default unicast forwarding classes.

- `best-effort`—Best-effort traffic
- `fcoe`—Guaranteed delivery for Fibre Channel over Ethernet traffic (do not use on OCX Series switches)
- `no-loss`—Guaranteed delivery for TCP no-loss traffic (do not use on OCX Series switches)
- `network-control`—Network control traffic

NOTE: QFX10002-60C does not support PFC and lossless queues; that is, default lossless queues (`fcoe` and `no-loss`) will be lossy queues.

The default multdestination forwarding class, except on QFX10000 switches and NFX Series devices, is:

- `mcast`—Multdestination traffic

The NFX Series devices have the following default forwarding classes:

- best-effort (be)—Provides no service profile. Loss priority is typically not carried in a CoS value.
- expedited-forwarding (ef)—Provides a low loss, low latency, low jitter, assured bandwidth, end-to-end service.
- assured-forwarding (af)—Provides a group of values you can define and includes four subclasses: AF1, AF2, AF3, and AF4, each with two drop probabilities: low and high.
- network-control (nc)—Supports protocol control and thus is typically high priority.

You can map forwarding classes to queues using the `class` statement. You can map more than one forwarding class to a single queue. Except on QFX10000 or NFX Series devices, all forwarding classes mapped to a particular queue must be of the same type, either unicast or multicast. You cannot mix unicast and multicast forwarding classes on the same queue.

All of the forwarding classes mapped to the same queue must have the same packet drop attribute: either all of the forwarding classes must be lossy or all of the forwarding classes must be lossless. This is important because the default fcoe and no-loss forwarding classes have the `no-loss` drop attribute, which is not supported on OCX Series switches. On OCX Series switches, do not map traffic to the default fcoe and no-loss forwarding classes.

```
[edit class-of-service forwarding-classes]
user@switch# set class class-name queue-num queue-number <no-loss>
```

One example is to create a forwarding class named `be2` and map it to queue 1:

```
[edit class-of-service forwarding-classes]
user@switch# set class be2 queue-num 1
```

Another example is to create a lossless forwarding class named `fcoe2` and map it to queue 5:

```
[edit class-of-service forwarding-classes]
user@switch# set class fcoe2 queue-num 5 no-loss
```

NOTE: On switches that do not run ELS software, if you are using Junos OS Release 12.2 or later, use the default forwarding-class-to-queue mapping for the lossless `fcoe` and `no-loss` forwarding classes. If you explicitly configure the lossless forwarding classes, the traffic mapped to those forwarding classes is treated as lossy (best-effort) traffic and does *not* receive lossless treatment

unless you include the optional `no-loss` packet drop attribute introduced in Junos OS Release 12.3 in the forwarding class configuration..

NOTE: On switches that do not run ELS software, Junos OS Release 11.3R1 and earlier supported an alternate method of mapping forwarding classes to queues that allowed you to map only one forwarding class to a queue using the statement:

```
[edit class-of-service forwarding-classes]
user@switch# set queue queue-number class-name
```

The `queue` statement has been deprecated and is no longer valid in Junos OS Release 11.3R2 and later. If you have a configuration that uses the `queue` statement to map forwarding classes to queues, edit the configuration to replace the `queue` statement with the `class` statement.

Release History Table

Release	Description
22.1R1	Starting in Junos OS Release 22.1R1, QFX10000 Series devices support 16 forwarding classes.

RELATED DOCUMENTATION

- Example: Configuring CoS Hierarchical Port Scheduling (ETS)*
- Example: Configuring Forwarding Classes*
- Monitoring CoS Forwarding Classes*
- Understanding CoS Forwarding Classes*
- Understanding CoS Port Schedulers on QFX Switches*

Example: Configuring Forwarding Classes

IN THIS SECTION

- Requirements | 134

- [Overview | 134](#)
- [Example 1: Configuring Forwarding Classes for Switches Except QFX10000 | 136](#)
- [Example 2: Configuring Forwarding Classes for QFX10000 Switches | 138](#)

Forwarding classes group packets for transmission. Forwarding classes map to output queues, so the packets assigned to a forwarding class use the output queue mapped to that forwarding class. Except on QFX10000, unicast traffic and multdestination (multicast, broadcast, and destination lookup fail) traffic use separate forwarding classes and output queues.

Requirements

This example uses the following hardware and software components for two configuration examples:

Configuring forwarding classes for switches except QFX10000

- One switch except QFX10000 (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

Configuring forwarding classes for QFX10000 switches

- One QFX10000 switch
- Junos OS Release 15.1X53-D10 or later for the QFX Series

Overview

The QFX10000 switch supports eight forwarding classes. Other switches support up to 12 forwarding classes. To forward traffic, you must map (assign) the forwarding classes to output queues. On the QFX10000 switch, queues 0 through 7 are for both unicast and multdestination traffic. On other switches, queues 0 through 7 are for unicast traffic, and queues 8 through 9 (QFX5200 switch) or 8 through 11 (other switches) are for multdestination traffic. Except for OCX Series switches, switches support up to six lossless forwarding classes. (OCX Series switches do not support lossless Layer 2 transport.)

The switch provides four default forwarding classes, and except on QFX10000 switches, these four forwarding classes are unicast, plus one default multdestination forwarding class. You can define the remaining forwarding classes and configure them as unicast or multdestination forwarding classes by mapping them to unicast or multdestination queues. The type of queue, unicast or multdestination, determines the type of forwarding class.

The four default forwarding classes (unicast except on QFX10000) are:

- `be`—Best-effort traffic
- `fcoe`—Guaranteed delivery for Fibre Channel over Ethernet traffic (do not use on OCX Series switches)
- `no-loss`—Guaranteed delivery for TCP no-loss traffic (do not use on OCX Series switches)
- `nc`—Network control traffic

Except on QFX10000 switches, the default multidestination forwarding class is:

- `mcast`—Multidestination traffic

Map forwarding classes to queues using the `class` statement. You can map more than one forwarding class to a single queue, but all forwarding classes mapped to a particular queue must be of the same type:

- Except on QFX10000 switches, all forwarding classes mapped to a particular queue must be either unicast or multicast. You cannot mix unicast and multicast forwarding classes on the same queue.
- On QFX10000 switches, all forwarding classes mapped to a particular queue must have the same packet drop attribute: all of the forwarding classes must be lossy, or all of the forwarding classes mapped to a queue must be lossless.

```
[edit class-of-service forwarding-classes]
user@switch# set class class-name queue-num queue-number;
```

NOTE: On switches that do not run ELS software, if you are using Junos OS Release 12.2, use the default forwarding-class-to-queue mapping for the lossless `fcoe` and `no-loss` forwarding classes. If you explicitly configure the lossless forwarding classes, the traffic mapped to those forwarding classes is treated as lossy (best-effort) traffic and does *not* receive lossless treatment. In Junos OS Release 12.3 and later, you can include the *no-loss* packet drop attribute in explicit forwarding class configurations to configure a lossless forwarding class.

NOTE: On switches that do not run ELS software, Junos OS Release 11.3R1 and earlier supported an alternate method of mapping forwarding classes to queues that allowed you to map only one forwarding class to a queue using the statement:

```
[edit class-of-service forwarding-classes]
user@switch# set queue queue-number class-name
```

The `queue` statement has been deprecated and is no longer valid in Junos OS Release 11.3R2 and later. If you have a configuration that uses the `queue` statement to map forwarding classes to queues, edit the configuration to replace the `queue` statement with the `class` statement.

NOTE: Hierarchical scheduling controls output queue forwarding. When you define a forwarding class and classify traffic into it, you must also define a scheduling policy for the forwarding class. Defining a scheduling policy means:

- Mapping a scheduler to the forwarding class in a scheduler map
- Including the forwarding class in a forwarding class set
- Associating the scheduler map with a traffic control profile
- Attaching the traffic control profile to a forwarding class set and applying the traffic control profile to an interface

On QFX10000 switches, you can define a scheduling policy using port scheduling:

- Mapping a scheduler to the forwarding class in a scheduler map.
- Applying the scheduler map to one or more interfaces.

Example 1: Configuring Forwarding Classes for Switches Except QFX10000

IN THIS SECTION

- [Verification](#) | 137

Configuration

Step-by-Step Procedure

[Table 43 on page 137](#) shows the configuration forwarding-class-to-queue mapping for this example:

Table 43: Forwarding-Class-to-Queue Example Configuration Except on QFX10000

Forwarding Class	Queue
best-effort	0
nc	7
mcast	8

To configure CoS forwarding classes for switches except QFX10000:

1. Map the best-effort forwarding class to queue 0:

```
[edit class-of-service forwarding-classes]
user@switch# set class best-effort queue-num 0
```

2. Map the nc forwarding class to queue 7:

```
[edit class-of-service forwarding-classes]
user@switch# set class nc queue-num 7
```

3. Map the mcast-be forwarding class to queue 8:

```
[edit class-of-service forwarding-classes]
user@switch# set class mcast-be queue-num 8
```

Verification

IN THIS SECTION

- [Verifying the Forwarding-Class-to-Queue Mapping | 138](#)

Verifying the Forwarding-Class-to-Queue Mapping

Purpose

Verify the forwarding-class-to-queue mapping. (The system shows only the explicitly configured forwarding classes; it does not show default forwarding classes such as fcoe and no-loss.)

Action

Verify the results of the forwarding class configuration using the operational mode command `show configuration class-of-service forwarding-classes`:

```
user@switch> show configuration class-of-service forwarding-classes
class best-effort queue-num 0;
class network-control queue-num 7;
class mcast queue-num 8;
```

Example 2: Configuring Forwarding Classes for QFX10000 Switches

IN THIS SECTION

Verification | 139

Configuration

Step-by-Step Procedure

[Table 44 on page 138](#) shows the configuration forwarding-class-to-queue mapping for this example:

Table 44: Forwarding-Class-to-Queue Example Configuration on QFX10000

Forwarding Class	Queue
best-effort	0
be1	1

Table 44: Forwarding-Class-to-Queue Example Configuration on QFX10000 (Continued)

Forwarding Class	Queue
nc	7

To configure CoS forwarding classes for QFX10000 switches:

1. Map the best-effort forwarding class to queue 0:

```
[edit class-of-service forwarding-classes]
user@switch# set class best-effort queue-num 0
```

2. Map the be1 forwarding class to queue 1:

```
[edit class-of-service forwarding-classes]
user@switch# set class be1 queue-num 1
```

3. Map the nc forwarding class to queue 7:

```
[edit class-of-service forwarding-classes]
user@switch# set class nc queue-num 7
```

Verification

IN THIS SECTION

- [Verifying the Forwarding-Class-to-Queue Mapping](#) | 140

Verifying the Forwarding-Class-to-Queue Mapping

Purpose

Verify the forwarding-class-to-queue mapping. (The system shows only the explicitly configured forwarding classes; it does not show default forwarding classes such as `fcoe` and `no-loss`.)

Action

Verify the results of the forwarding class configuration using the operational mode command `show configuration class-of-service forwarding-classes`:

```
user@switch> show configuration class-of-service forwarding-classes
class best-effort queue-num 0;
class be1 queue-num 1;
class network-control queue-num 7;
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Defining CoS Forwarding Classes

Monitoring CoS Forwarding Classes

[Overview of CoS Changes Introduced in Junos OS Release 11.3 | 57](#)

[Overview of CoS Changes Introduced in Junos OS Release 12.2 | 67](#)

Understanding CoS Forwarding Classes

[Understanding CoS Forwarding Classes](#)

Understanding CoS Forwarding Class Sets (Priority Groups)

A forwarding class set is the Junos OS configuration construct that equates to a priority group in enhanced transmission selection (ETS, described in IEEE 802.1Qaz). The switch implements ETS using a two-tier hierarchical scheduler.

A priority group is a group of forwarding classes. Each forwarding class is mapped to an output queue and an IEEE 802.1p priority (code points). Classifying traffic into a forwarding class based on its code points, and mapping the forwarding class to a queue, defines the traffic assigned to that queue. The forwarding classes that belong to a priority group share the port bandwidth allocated to that priority

group. The traffic mapped to forwarding classes in one priority group usually shares similar traffic-handling requirements.

You can configure up to three unicast forwarding class sets and one multicast forwarding class set. Only unicast forwarding classes can belong to unicast forwarding class sets. Only multicast forwarding classes can belong to the multicast forwarding class set.

If you configure a strict-high priority forwarding class (you can configure only one strict-high priority forwarding class), you must observe the following rules when configuring forwarding class sets:

- You must create a separate forwarding class set for the strict-high priority forwarding class.
- Only one forwarding class set can contain the strict-high priority forwarding class.
- A strict-high priority forwarding class cannot belong to the same forwarding class set as forwarding classes that are not strict-high priority.
- A strict-high priority forwarding class cannot belong to a multidestination forwarding class set.
- You cannot configure a guaranteed minimum bandwidth (guaranteed rate) for a forwarding class set that includes a strict-high priority forwarding class. (You also cannot configure a guaranteed minimum bandwidth for a strict-high forwarding class.)
- We recommend that you always apply a shaping rate to a strict-high priority forwarding class to prevent it from starving the queues mapped to other forwarding classes. If you do not apply a shaping rate to limit the amount of bandwidth a strict-high priority forwarding class can use, then the strict-high priority forwarding class can use all of the available port bandwidth and starve other forwarding classes on the port.

You must use hierarchical scheduling if you explicitly configure CoS. The two-tier hierarchical scheduler defines bandwidth resources for the forwarding class set (priority group), and then allocates those resources among the forwarding classes (priorities) that belong to the forwarding class set.

If you do not explicitly configure forwarding class sets, the system automatically creates a default forwarding class set that contains all of the forwarding classes on the switch. The system assigns 100 percent of the port output bandwidth to the default forwarding class set. Ingress traffic is classified based on the default classifier settings. The forwarding classes in the default forwarding class set receive bandwidth based on the default scheduler settings. Forwarding classes that are not part of the default scheduler receive no bandwidth. The default priority group is transparent. It does not appear in the configuration and is used for Data Center Bridging Capability Exchange Protocol (DCBX) advertisement (except on OCX Series switches, which do not support DCBX).

When you explicitly configure forwarding class sets and apply them to interfaces, on those interfaces, forwarding classes that you do not map to a forwarding class set receive no guaranteed bandwidth. Forwarding classes that belong to the default forwarding class set might receive bandwidth if the other forwarding class sets are not using all of the port bandwidth. However, the amount of bandwidth received by forwarding classes that are not members of a forwarding class set is not guaranteed. In this

case, the bandwidth a forwarding class receives if it is not a member of a forwarding class set depends on whether unused port bandwidth is available and therefore is not deterministic.

To guarantee bandwidth for forwarding classes in a predictable manner, be sure to map all forwarding classes that you expect to carry traffic on an interface to a forwarding class set, and apply the forwarding class set to the interface.

RELATED DOCUMENTATION

Understanding CoS Hierarchical Port Scheduling (ETS)

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Forwarding Class Sets

Defining CoS Forwarding Class Sets

Defining CoS Forwarding Class Sets

A forwarding class set is a priority group for enhanced transmission selection (ETS) traffic control. Each forwarding class set consists of one or more forwarding classes. Classifiers map traffic into forwarding classes based on code points (priority), and forwarding classes are mapped to output queues.

You can configure up to three unicast forwarding class sets and one multicast forwarding class set.

To configure a forwarding class set using the CLI:

1. Assign one or more forwarding classes to the forwarding class set:

```
[edit class-of-service]
user@switch# set forwarding-class-sets forwarding-class-set-name class forwarding-class-name
```

2. Map the forwarding class set to an interface:

```
[edit class-of-service]
user@switch# set interfaces interface-name forwarding-class-set forwarding-class-set-name
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Forwarding Class Sets

Defining CoS Queue Schedulers

Defining CoS Traffic Control Profiles (Priority Group Scheduling)

Understanding CoS Forwarding Class Sets (Priority Groups)

Example: Configuring Forwarding Class Sets

IN THIS SECTION

- [Requirements | 144](#)
- [Overview | 144](#)
- [Verification | 146](#)

A forwarding class set (fc-set) is a priority group for enhanced transmission selection (ETS) traffic control. Each fc-set consists of one or more forwarding classes (priorities). Classifiers map traffic to forwarding classes based on code points, and forwarding classes are mapped to output queues.

ETS enables you to configure link resources (bandwidth and bandwidth sharing characteristics) for an fc-set, and then allocate the fc-set's resources among the forwarding classes that belong to the fc-set. This is called two-tier, or hierarchical, scheduling. Traffic control profiles control the scheduling for the fc-set (priority group), and schedulers control the scheduling for individual forwarding classes (priorities).

Configuring Forwarding Class Sets

Step-by-Step Procedure

1. Define the lan-pg priority group (fc-set) and assign to it the forwarding classes best-effort-1 and best-effort-2:

```
[edit class-of-service]
user@switch# set forwarding-class-sets lan-pg class best-effort-1
user@switch# set forwarding-class-sets lan-pg class best-effort-2
```

2. Define the san-pg priority group and assign to it the forwarding classes fcoe and fcoe-2:

```
[edit class-of-service]
user@switch# set forwarding-class-sets san-pg class fcoe
user@switch# set forwarding-class-sets san-pg class fcoe-2
```

3. Define the hpc-pg priority group and assign to it the forwarding classes nc and high-perf:

```
[edit class-of-service]
user@switch# set forwarding-class-sets hpc-pg class nc
user@switch# set forwarding-class-sets hpc-pg class high-perf
```

4. Map the three forwarding class sets to an interface (the output traffic control profiles associated with the forwarding class sets determine the class of service scheduling for the priority groups):

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/7 forwarding-class-set lan-pg output-traffic-control-
profile lan-tcp
user@switch# set interfaces xe-0/0/7 forwarding-class-set san-pg output-traffic-control-
profile san-tcp
user@switch# set interfaces xe-0/0/7 forwarding-class-set hpc-pg output-traffic-control-
profile hpc-tcp
```

Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series.

Overview

You can configure up to three unicast fc-sets and one multicast fc-set. A common way to configure unicast priority groups is to configure separate fc-sets for local area network (LAN) traffic, storage area network (SAN) traffic, and high-performance computing (HPC) traffic, and then assign the appropriate forwarding classes to each fc-set.

NOTE: If you configure a strict-high priority forwarding class, you must create an fc-set that is dedicated only to strict-high priority traffic. You can only configure one strict-high priority forwarding class, and only one fc-set can contain a strict-high priority queue. Queues that are not strict-high priority cannot belong to the same fc-set as a strict-high priority queue. The multidestination fc-set cannot contain a strict-high priority queue.

To apply ETS, you use a traffic control profile to map one or more fc-sets to a physical egress port. You can map up to three unicast forwarding class sets and one multidestination forwarding class set to each port. When you map an fc-set to a port, the port uses hierarchical scheduling to allocate port resources to the priority group (fc-set) and to allocate the priority group resources to the queues (forwarding classes) that belong to the priority group.

This example describes how to:

- Configure three fc-sets called lan-pg, san-pg, and hpc-pg.
- Assign forwarding classes to each of the fc-sets.
- Apply the fc-sets and their output traffic control profiles to an egress interface.

This example does not describe how to configure the forwarding classes assigned to the fc-sets or how to configure traffic control profiles (scheduling). [Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#) provides a complete example of how to configure ETS, including forwarding class and scheduling configuration. [Table 45 on page 145](#) shows the configuration components for this example:

Table 45: Components of the Forwarding Class Sets Configuration Example

Component	Settings
Hardware	QFX3500 switch
LAN traffic priority group	Forwarding class set: lan-pg Forwarding classes: best-effort-1, best-effort-2

Table 45: Components of the Forwarding Class Sets Configuration Example *(Continued)*

Component	Settings
SAN traffic priority group	Forwarding class set: san-pg Forwarding classes: fcoe, fcoe-2 NOTE: OCX Series switches do not support FCoE traffic or lossless Layer 2 transport. If you were configuring this example on an OCX Series switch, you could omit this priority group, or rename it and map different forwarding classes to it.
HPC traffic priority group	Forwarding class set: hpc-pg Forwarding classes: nc, high-perf
Egress interface	xe-0/0/7

Verification

IN THIS SECTION

- [Verifying Forwarding Class Set Membership | 146](#)
- [Verifying the Egress Interface Configuration | 147](#)

To verify the priority group configuration, perform these tasks:

Verifying Forwarding Class Set Membership

Purpose

Verify that you configured the lan-pg, san-pg, and hpc-pg priority groups with the correct forwarding classes.

Action

List the forwarding class set member configuration using the operational mode command `show configuration class-of-service forwarding-class-sets`:

```
user@switch> show configuration class-of-service forwarding-class-sets
lan-pg {
    class best-effort-1;
    class best-effort-2;
}
san-pg {
    class fcoe;
    class fcoe-2;
}
hpc-pg {
    class high-perf;
    class nc;
}
```

Verifying the Egress Interface Configuration

Purpose

Verify that egress interface `xe-0/0/7` is associated with the `lan-pg`, `san-pg`, and `hpc-pg` priority groups and with the correct output traffic control profiles.

Action

Display the egress interface using the operational mode command `show configuration class-of-service interfaces xe-0/0/7`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/7
forwarding-class-set {
    lan-pg {
        output-traffic-control-profile lan-tcp;
    }
    san-pg {
        output-traffic-control-profile san-tcp;
    }
    hpc-pg {
```

```

        output-traffic-control-profile hpc-tcp;
    }
}

```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Queue Schedulers

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Defining CoS Forwarding Class Sets

Understanding CoS Forwarding Class Sets (Priority Groups)

Understanding Host Routing Engine Outbound Traffic Queues and Defaults

The host Routing Engine and CPU generate outbound traffic that is transmitted using different protocols. You cannot configure a classifier to map different types of outbound traffic that the host generates to forwarding classes (queues). The traffic that the host generates is assigned to forwarding classes by default as shown in [Table 46 on page 149](#).

If you want to separate host outbound traffic from other traffic or if you want to assign that traffic to a particular queue, you can configure a single forwarding class for all traffic that the host generates. If you configure a forwarding class for outbound host traffic, that forwarding class is used globally for all traffic generated by the host. (That is, the host outbound traffic is mapped to the selected queue on all egress interfaces.) Configuring a forwarding class for host outbound traffic does not affect transit or incoming traffic.

Whether you use the default host outbound traffic forwarding class configuration or configure a forwarding class for all host outbound traffic, the configuration applies to all Layer 2 and Layer 3 protocols and to all application-level traffic such as FTP and ping operations.

If you configure a queue for host outbound traffic, the queue must be properly configured on all interfaces.

NOTE: Fibre Channel over Ethernet (FCoE) Initialization Protocol (FIP) packets generated by the CPU are always transmitted on the `fcoe` queue (queue 3), even if you configure a queue for host outbound traffic. This helps to ensure lossless behavior for FCoE traffic. QFabric systems classify

FIP control packets into the same traffic class (fcoe) across the Interconnect device (fabric) and the egress Node device.

This does not apply to OCX Series switches, which do not support FCoE.

By default, traffic generated by the host is sent to the best effort queue (queue 0) or to the network control queue (queue 7). [Table 46 on page 149](#) lists the default host traffic to output queue mapping.

Table 46: Routing Engine Protocol Default Queue Mapping

Routing Engine Protocol	Default Queue Mapping
Address Resolution Protocol (ARP) reply	Queue 0
ARP request	Queue 0
Bidirectional Forwarding Detection (BFD) Protocol	Queue 7
Border Gateway Protocol (BGP)	Queue 0
BGP TCP Retransmission	Queue 7
Fibre Channel over Ethernet (FCoE) Initialization Protocol (FIP)	Queue 3
File Transfer Protocol (FTP)	Queue 0
Internet Control Message Protocol (ICMP) reply	Queue 0
ICMP request	Queue 0
Internet Group Management Protocol (IGMP) query	Queue 7
IGMP report	Queue 0
Link Aggregation Control Protocol (LACP)	Queue 7

Table 46: Routing Engine Protocol Default Queue Mapping *(Continued)*

Routing Engine Protocol	Default Queue Mapping
Open Shortest Path First (OSPF) hello	Queue 7
OSPF protocol data unit (PDU)	Queue 7
OSPF link state advertisements (LSAs)	Queue 7
Protocol Independent Multicast (PIM)	Queue 7
PIM hello	Queue 7
Simple Network Management Protocol (SNMP)	Queue 0
Secure Shell (SSH)	Queue 0
Telnet	Queue 0
Virtual Router Redundancy Protocol (VRRP)	Queue 7
VLAN Spanning Tree Protocol (VSTP)	Queue 7
xnm-clear-text	Queue 0
xnm-ssl	Queue 0

RELATED DOCUMENTATION

Understanding CoS Forwarding Classes

[Understanding CoS Forwarding Classes](#)

Changing the Host Outbound Traffic Default Queue Mapping

Example: Configuring Forwarding Classes

Changing the Host Outbound Traffic Default Queue Mapping

If you do not want to use the default mapping of host Routing Engine and CPU outbound traffic to queues, you can change the default output queue. You can also change the default DSCP bits used in the type of service (ToS) field of packets generated by the Routing Engine.

Configuring a queue for host outbound traffic maps all traffic that the host generates to one forwarding class (queue). The configuration is global and applies to all host-generated traffic on the switch. Configuring a forwarding class for host outbound traffic does not affect transit or incoming traffic.

NOTE: Fibre Channel over Ethernet (FCoE) Initialization Protocol (FIP) packets generated by the CPU are always transmitted on the `fcoe` queue (queue 3), even if you configure a queue for host outbound traffic. This helps to ensure lossless behavior for FCoE traffic. QFabric systems classify FIP control packets into the same traffic class (`fcoe`) across the Interconnect device (fabric) and the egress Node device.

This does not apply to OCX Series switches, which do not support FCoE.

To change the host outbound traffic egress queue by including the `host-outbound-traffic` statement at the `[edit class-of-service]` hierarchy level:

```
[edit class-of-service]
host-outbound-traffic {
    forwarding-class class-name;
    dscp-code-point code-point;
}
```

For example, to map host outbound traffic to queue 7 (the network control forwarding class) and set the DSCP code point value to 101010:

```
[edit class-of-service]
host-outbound-traffic {
    forwarding-class network-control;
    dscp-code-point 101010
}
```

RELATED DOCUMENTATION

Understanding Host Routing Engine Outbound Traffic Queues and Defaults

Understanding CoS Rewrite Rules

As packets enter or exit a network, edge switches might be required to alter the class-of-service (CoS) settings of the packets. *Rewrite rules* set the value of the code point bits (Layer 3 DSCP bits, Layer 2 CoS bits, or MPLS EXP bits) within the header of the outgoing packet. Each rewrite rule:

1. Reads the current forwarding class and loss priority associated with the packet.
2. Locates the new (rewrite) code point value from a table.
3. Writes that code point value into the packet header, replacing the old code point value.

Rewrite rules must be assigned to an interface for rewrites to take effect.

You can apply (bind) one DSCP or DSCP IPv6 rewrite rule and one IEEE 802.1p rewrite rule to each interface. You can also bind EXP rewrite rules to `family mpls` logical interfaces to rewrite the CoS bits of MPLS traffic.

NOTE: OCX Series switches do not support MPLS and do not support EXP rewrite rules.

You cannot apply both a DSCP and a DSCP IPv6 rewrite rule to the same physical interface. Each physical interface supports only one DSCP rewrite rule. Both IP and IPv6 packets use the same DSCP rewrite rule, regardless if the configured rewrite rule is DSCP or DSCP IPv6. You can apply an EXP rewrite rule on an interface that has DSCP or IEEE rewrite rules. Only MPLS traffic on `family mpls` interfaces uses the EXP rewrite rule.

You *can* apply both a DSCP rewrite rule and a DSCP IPv6 rewrite rule to a logical interface. IPv6 packets are rewritten with DSCP-IPv6 rewrite-rules and IPv4 packets are remarked with DSCP rewrite-rules.

NOTE: There are no default rewrite rules. If you want to apply a rewrite rule to outgoing packets, you must explicitly configure the rewrite rule.

You can look at behavior aggregate (BA) classifiers and rewrite rules as two sides of the same coin. A BA classifier reads the code point bits of incoming packets and classifies the packets into forwarding classes, then the system applies the CoS configured for the forwarding class to those packets. Rewrite rules change (rewrite) the code point bits just before the packets leave the system so that the next switch or router can apply the appropriate level of CoS to the packets. When you apply a rewrite rule to an interface, the rewrite rule is the last CoS action performed on the packet before it is forwarded.

Rewrite rules alter CoS values in outgoing packets on the outbound interfaces of an edge switch to accommodate the policies of a targeted peer. This allows the downstream switch in a neighboring network to classify each packet into the appropriate service group.

NOTE: On each physical interface, either all forwarding classes that are being used on the interface must have rewrite rules configured or no forwarding classes that are being used on the interface can have rewrite rules configured. On any physical port, do not mix forwarding classes with rewrite rules and forwarding classes without rewrite rules.

NOTE: Rewrite rules are applied *before* the egress filter is matched to traffic. Because the code point rewrite occurs before the egress filter is matched to traffic, the egress filter match is based on the rewrite value, not on the original code point value in the packet.

For packets that carry both an inner VLAN tag and an outer VLAN tag, the rewrite rule rewrites only the outer VLAN tag.

MPLS EXP rewrite rules apply only to family `mpls` logical interfaces. You cannot apply to an EXP rewrite rule to a physical interface. You can configure up to 64 EXP rewrite rules, but you can only use 16 EXP rewrite rules at any time on the switch. On a given logical interface, all pushed MPLS labels have the same EXP rewrite rule applied to them. You can apply different EXP rewrite rules to different logical interfaces on the same physical interface.

NOTE: If the switch is performing penultimate hop popping (PHP), EXP rewrite rules do not take effect. If both an EXP classifier and an EXP rewrite rule are configured on the switch, then the EXP value from the last popped label is copied into the inner label. If either an EXP classifier or an EXP rewrite rule (but not both) is configured on the switch, then the inner label EXP value is sent unchanged.

You can configure enough rewrite rules to handle most, if not all, network scenarios. [Table 47 on page 153](#) shows how many of each type of rewrite rules you can configure, and how many entries you can configure per rewrite rule.

Table 47: Configuring Rewrite Rules

Rewrite Rule Type	Maximum Number of Rewrite Rules	Maximum Number of Entries per Rewrite Rule
IEEE 802.1p	64	128
DSCP	32	128

Table 47: Configuring Rewrite Rules *(Continued)*

Rewrite Rule Type	Maximum Number of Rewrite Rules	Maximum Number of Entries per Rewrite Rule
DSCP IPv6	32	128
MPLS EXP	64	128

You cannot apply rewrite rules directly to integrated routing and bridging (IRB), also known as routed VLAN interfaces (RVIs), because the members of IRBs/RVIs are VLANs, not ports. However, you can apply rewrite rules to the VLAN port members of an IRB/*RVI*.

NOTE: OCX Series switches do not support IRBs/RVIs.

RELATED DOCUMENTATION

Understanding Junos CoS Components

Defining CoS Rewrite Rules

Configuring Rewrite Rules for MPLS EXP Classifiers

Defining CoS Rewrite Rules

Overview

Edge switches might need to change the class-of-service (CoS) settings of the packets. You can configure rewrite rules to alter code point bit values in outgoing packets on the outbound interfaces of a switch so that the CoS treatment matches the policies of a targeted peer. Policy matching allows the downstream routing platform or switch in a neighboring network to classify each packet into the appropriate service group.

To configure a CoS rewrite rule, create the rule by giving it a name and associating it with a forwarding class, loss priority, and code point. This creates a rewrite table. After the rewrite rule is created, enable it on an interface (EXP rewrite rules can only be enabled on `family mpls` logical interfaces, not on physical interfaces). You can also apply an existing rewrite rule on an interface.

NOTE: On each physical interface, either all forwarding classes that are being used on the interface must have rewrite rules configured, or no forwarding classes that are being used on the interface can have rewrite rules configured. On any physical port, do not mix forwarding classes with rewrite rules and forwarding classes without rewrite rules.

NOTE: To replace an existing rewrite rule on the interface with a new rewrite rule of the same type, first explicitly remove the existing rewrite rule and then apply the new rule.

NOTE: For packets that carry both an inner VLAN tag and an outer VLAN tag, the rewrite rule rewrites only the outer VLAN tag.

Platform-specific Information

- OCX Series switches do not support MPLS, so they do not support EXP rewrite rules.
- QFX5130, QFX5700 & QFX5220 switches do not support DSCP IPv6 classifiers and rewrite rules. However, you can apply DSCP classifiers and rewrite rules for IPV6 traffic as well.

Configuring Rewrite Rules

To create rewrite rules and enable them on interfaces:

- To create an 802.1p rewrite rule named `customup-rw` in the rewrite table for all Layer 2 interfaces:

```
[edit class-of-service rewrite-rules]
user@switch# set ieee-802.1 customup-rw forwarding-class be loss-priority low code-point 000
user@switch# set ieee-802.1 customup-rw forwarding-class be loss-priority high code-point 001
user@switch# set ieee-802.1 customup-rw forwarding-class be loss-priority low code-point 010
user@switch# set ieee-802.1 customup-rw forwarding-class fcoe loss-priority low code-point 011
user@switch# set ieee-802.1 customup-rw forwarding-class ef-no-loss loss-priority low code-point 100
user@switch# set ieee-802.1 customup-rw forwarding-class ef-no-loss loss-priority high code-point 101
user@switch# set ieee-802.1 customup-rw forwarding-class nc loss-priority low code-point 110
user@switch# set ieee-802.1 customup-rw forwarding-class nc loss-priority high code-point 111
```

- To enable an 802.1p rewrite rule named `customup-rw` on a Layer 2 interface:

```
[edit]
user@switch# set class-of-service interfaces xe-0/0/7 unit 0 rewrite-rules ieee-802.1
customup-rw
```

NOTE: All forwarding classes assigned to port `xe-0/0/7` must have rewrite rules. Do not mix forwarding classes that have rewrite rules with forwarding classes that do not have rewrite rules on the same physical interface.

- To enable an 802.1p rewrite rule named `customup-rw` on all 10-Gigabit Ethernet interfaces on the switch, use wildcards for the interface name and logical interface (unit) number:

```
[edit]
user@switch# set class-of-service interfaces xe-* unit * rewrite-rules customup-rw
```

NOTE: In this case, *all* forwarding classes assigned to *all* 10-Gigabit Ethernet ports must have rewrite rules. Do not mix forwarding classes that have rewrite rules with forwarding classes that do not have rewrite rules on the same physical interface.

RELATED DOCUMENTATION

Monitoring CoS Rewrite Rules

Configuring Rewrite Rules for MPLS EXP Classifiers

Understanding CoS Rewrite Rules

Understanding CoS MPLS EXP Classifiers and Rewrite Rules

Configuring CoS Fixed Classifier Rewrite Values for Native FC Interfaces (NP_Ports)

Fibre Channel over Ethernet (FCoE) traffic typically uses IEEE 802.1p priority 3 (code point 011). When Fibre Channel (FC) traffic arrives on a native FC interface (NP_Port) on an FCoE-FC gateway, the interface encapsulates the FC traffic in Ethernet to create FCoE frames. By default, the native FC

interface assigns priority 3 to the FCoE traffic. The traffic is then forwarded internally to the gateway Ethernet interfaces, and then forwarded to the FCoE network.

If your FCoE network uses priority 3 for FCoE traffic, you do not need to use a rewrite value to remap the FCoE priority on native FC interfaces, because the default configuration maps priority 3 to the FCoE forwarding class.

However, if the FCoE network uses a different priority than priority 3 for FCoE traffic, then you can configure a rewrite value to remap incoming traffic from the FC SAN to that priority after the interface encapsulates the FC packets in Ethernet. Setting a rewrite value for the IEEE 802.1p code point (priority) configures the gateway native FC interface to assign the rewrite value to the encapsulated FCoE frames before forwarding the FCoE frames to the gateway Ethernet interface. Instead of a priority of 3, the FCoE frames use the priority specified in the rewrite value.

Traffic coming from the FC SAN is classified into a lossless forwarding class, and that lossless forwarding class is mapped to the rewrite value (the priority used for FCoE traffic on the converged Ethernet network). You specify the lossless forwarding class used for FCoE traffic on a native FC interface by configuring a fixed classifier and applying it to the native FC interface. (The same forwarding class must also be mapped to the rewrite value priority in the ingress classifier applied to the FCoE Ethernet interfaces.) All traffic received from the FC SAN on that FC interface is encapsulated in Ethernet, classified into the forwarding class specified in the fixed classifier, and assigned the rewrite value priority.

Configuring a rewrite value consists of:

- Configuring a fixed classifier on the native FC interface. The fixed classifier assigns all the traffic that arrives at the interface from the connected peer in the FC SAN to one fixed forwarding class. The forwarding class must be a lossless forwarding class and must be classified to the rewrite value priority in the ingress classifier configuration on the FCoE Ethernet interfaces.
- Specifying an IEEE 802.1p rewrite value for the native FC interface. The traffic mapped to the forwarding class in the fixed classifier is marked with the priority you specify in the rewrite value when the traffic is encapsulated in Ethernet. The rewrite value must be the IEEE 802.1p priority used for FCoE traffic in your converged Ethernet network.

You can configure one rewrite value for each local FCoE-FC gateway fabric. All of the native FC interfaces in a particular fabric must use the same rewrite value. Native FC interfaces that belong to different FCoE-FC gateway fabrics can use different rewrite values.

1. Configure a fixed classifier on the native FC interface:

[edit [class-of-service](#)]

```
user@switch# set interfaces fc-interface-name forwarding-class lossless-forwarding-class-name
```


For example, to configure a fixed classifier on native FC interface `fc-0/0/2` that specifies the lossless forwarding class `fcoe1`:

```
[edit class-of-service]
user@switch# set interfaces fc-0/0/2 forwarding-class fcoe1
```

2. Configure a rewrite value for the traffic classified into the fixed classifier (this must be the IEEE 802.1p priority used for the traffic on your converged Ethernet network):

```
[edit class-of-service]
user@switch# set interfaces fc-interface-name rewrite-value input ieee-802.1 code-point code-point-bits
```

For example, to configure a rewrite value on native FC interface `fc-0/0/2` that specifies an IEEE 802.1p priority of 101 (the lossless forwarding class specified in the fixed classifier must be classified to this priority in the ingress classifier configuration on the FCoE Ethernet interfaces):

```
[edit class-of-service]
user@switch# set interfaces fc-0/0/2 rewrite-value input ieee-802.1 code-point 101
```

In the example, all traffic from the FC SAN that arrives at FCoE-FC gateway interface `fc-0/0/2` is encapsulated in Ethernet, classified into the lossless `fcoe1` forwarding class, and tagged with the IEEE 802.1p priority 5 (code point 101). In this example, we assume that the converged Ethernet network uses priority 5 for FCoE traffic, and that the `fcoe1` forwarding class is mapped to priority 5 in the ingress classifier configuration on the Ethernet interfaces. To achieve lossless transport, you must also enable PFC on priority 5 on the Ethernet interfaces that connect the FCoE traffic to the Ethernet network.

RELATED DOCUMENTATION

[Example: Configuring IEEE 802.1p Priority Remapping on an FCoE-FC Gateway | 624](#)

[Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway | 620](#)

Troubleshooting an Unexpected Rewrite Value

IN THIS SECTION

- Problem | 159
- Cause | 159
- Solution | 160

Problem

Description

Traffic from one or more forwarding classes on an egress port is assigned an unexpected rewrite value.

NOTE: For packets that carry both an inner VLAN tag and an outer VLAN tag, the rewrite rules rewrite only the outer VLAN tag.

Cause

If you configure a rewrite rule for a forwarding class on an egress port, but you do not configure a rewrite rule for every forwarding class on that egress port, then the forwarding classes that do not have a configured rewrite rule are assigned random rewrite values.

For example:

1. Configure forwarding classes fc1, fc2, and fc3.
2. Configure rewrite rules for forwarding classes fc1 and fc2, but not for forwarding class fc3.
3. Assign forwarding classes fc1, fc2, and fc3 to a port.

When traffic for these forwarding classes flows through the port, traffic for forwarding classes fc1 and fc2 is rewritten correctly. However, traffic for forwarding class fc3 is assigned a random rewrite value.

Solution

If any forwarding class on an egress port has a configured rewrite rule, then all forwarding classes on that egress port must have a configured rewrite rule. Configuring a rewrite rule for any forwarding class that is assigned a random rewrite value solves the problem.

TIP: If you want the forwarding class to use the same code point value assigned to it by the ingress classifier, specify that value as the rewrite rule value. For example, if a forwarding class has the IEEE 802.1 ingress classifier code point value 011, configure a rewrite rule for that forwarding class that uses the IEEE 802.1p code point value 011.

NOTE: There are no default rewrite rules. You can bind one rewrite rule for DSCP traffic and one rewrite rule for IEEE 802.1p traffic to an interface. A rewrite rule can contain multiple forwarding-class-to-rewrite-value mappings.

1. To assign a rewrite value to a forwarding class, add the new rewrite value to the same rewrite rule as the other forwarding classes on the port:

```
[edit class-of-service rewrite-rules]
user@switch# set (dscp | ieee-802.1) rewrite-name forwarding-class class-name loss-priority
priority code-point (alias | bits)
```

For example, if the other forwarding classes on the port use rewrite values defined in the rewrite rule *custom-rw*, the forwarding class *be2* is being randomly rewritten, and you want to use IEEE 802.1 code point 002 for the *be2* forwarding class:

```
[edit class-of-service rewrite-rules]
user@switch# set ieee-802.1 custom-rw forwarding-class be2 loss-priority low code-point 002
```

2. Enable the rewrite rule on an interface if it is not already enabled on the desired interface:

```
[edit]
user@switch# set class-of-service interfaces interface-name unit unit rewrite-rules (dscp |
ieee-802.1) rewrite-rule-name
```

For example, to enable the rewrite rule `custom-rw` on interface `xe-0/0/24.0`:

```
[edit]
user@switch# set class-of-service interfaces xe-0/0/24 unit 0 rewrite-rules ieee-802.1 custom-
rw
```

RELATED DOCUMENTATION

[interfaces](#)

[rewrite-rules](#)

[Defining CoS Rewrite Rules](#)

[Monitoring CoS Rewrite Rules](#)

3

PART

Scheduling Traffic

[Using Schedulers \(Node Devices\) | 163](#)

[Using Schedulers \(Interconnect Device Fabric\) | 315](#)

Using Schedulers (Node Devices)

IN THIS CHAPTER

- Understanding Default CoS Scheduling and Classification | 164
- Understanding CoS Scheduling on QFabric System Node Device Fabric (fte) Ports | 174
- Understanding CoS Scheduling Behavior and Configuration Considerations | 179
- Understanding CoS Output Queue Schedulers | 186
- Defining CoS Queue Schedulers | 194
- Example: Configuring Queue Schedulers | 198
- Defining CoS Queue Scheduling Priority | 206
- Example: Configuring Queue Scheduling Priority | 207
- Understanding CoS Traffic Control Profiles | 213
- Understanding CoS Priority Group Scheduling | 214
- Defining CoS Traffic Control Profiles (Priority Group Scheduling) | 218
- Example: Configuring Traffic Control Profiles (Priority Group Scheduling) | 220
- Understanding CoS Hierarchical Port Scheduling (ETS) | 223
- Example: Configuring CoS Hierarchical Port Scheduling (ETS) | 230
- Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth | 266
- Understanding CoS Priority Group Shaping and Queue Shaping (Maximum Bandwidth) | 269
- Example: Configuring Minimum Guaranteed Output Bandwidth | 272
- Troubleshooting Egress Bandwidth That Exceeds the Configured Minimum Bandwidth | 279
- Example: Configuring Maximum Output Bandwidth | 280
- Troubleshooting Egress Bandwidth That Exceeds the Configured Maximum Bandwidth | 287
- Troubleshooting Egress Queue Bandwidth Impacted by Congestion | 288
- Understanding CoS WRED Drop Profiles | 290
- Configuring CoS WRED Drop Profiles | 297
- Example: Configuring WRED Drop Profiles | 300
- Configuring CoS Drop Profile Maps | 307
- Example: Configuring Drop Profile Maps | 307

- [Troubleshooting a Port Reset on QFabric Systems When a Queue Stops Transmitting Traffic | 311](#)

Understanding Default CoS Scheduling and Classification

IN THIS SECTION

- [Default Classification | 165](#)
- [Default Scheduling | 170](#)
- [Default DCBX Advertisement | 173](#)
- [Default Scheduling and Classification Summary | 174](#)

If you do not explicitly configure classifiers and apply them to interfaces, the switch uses the default classifier to group ingress traffic into forwarding classes. If you do not configure scheduling on an interface, the switch uses the default schedulers to provide egress port resources for traffic. Default classification maps all traffic into default forwarding classes (best-effort, fcoe, no-loss, network-control, and mcast). Each default forwarding class has a default scheduler, so that the traffic mapped to each default forwarding class receives port bandwidth, prioritization, and packet drop characteristics.

The switch supports direct port scheduling and enhanced transmission selection (ETS), also known as hierarchical port scheduling, except on QFX5200 and QFX5210 switches.

Hierarchical scheduling groups IEEE 802.1p priorities (IEEE 802.1p code points, which classifiers map to forwarding classes, which in turn are mapped to output queues) into priority groups (forwarding class sets). If you use only the default traffic scheduling and classification, the switch automatically creates a default priority group that contains all of the priorities (which are mapped to forwarding classes and output queues), and assigns 100 percent of the port output bandwidth to that priority group. The forwarding classes (queues) in the default forwarding class set receive bandwidth based on the default classifier settings. The default priority group is transparent. It does not appear in the configuration and is used for Data Center Bridging Capability Exchange (DCBX) protocol advertisement.

NOTE: If you explicitly configure one or more priority groups on an interface, any forwarding class that is not assigned to a priority group on that interface receives *no bandwidth*. This means that if you configure hierarchical scheduling on an interface, every forwarding class (priority) that

you want to forward traffic on that interface must belong to a forwarding class set (priority group). ETS is not supported on QFX5200 or QFX5210 switches.

The following sections describe:

Default Classification

On switches except QFX10000 and NFX Series devices, the default classifiers assign unicast and multicast best-effort and network-control ingress traffic to default forwarding classes and loss priorities. The switch applies default unicast IEEE 802.1, unicast DSCP, and multidestination classifiers to each interface that does not have explicitly configured classifiers.

On QFX10000 switches and NFX Series devices, the default classifiers assign ingress traffic to default forwarding classes and loss priorities. The switch applies default IEEE 802.1, DSCP, and DSCP IPv6 classifiers to each interface that does not have explicitly configured classifiers. If you do not configure and apply EXP classifiers for MPLS traffic to logical interfaces, MPLS traffic on interfaces configured as family mpls uses the IEEE classifier.

If you explicitly configure one type of classifier but not other types of classifiers, the system uses only the configured classifier and does not use default classifiers for other types of traffic. There are two default IEEE 802.1 classifiers: a trusted classifier for ports that are in trunk mode or tagged-access mode, and an untrusted classifier for ports that are in access mode.

NOTE: The default classifiers apply to unicast traffic except on QFX10000 switches and NFX Series devices. Tagged-access mode does not apply to QFX10000 switches or NFX Series devices.

[Table 48 on page 165](#) shows the default mapping of IEEE 802.1 code-point values to forwarding classes and loss priorities for ports in trunk mode or tagged-access mode.

Table 48: Default IEEE 802.1 Classifiers for Ports in Trunk Mode or Tagged-Access Mode (Trusted Classifier)

Code Point	Forwarding Class	Loss Priority
be (000)	best-effort	low
be1 (001)	best-effort	low

Table 48: Default IEEE 802.1 Classifiers for Ports in Trunk Mode or Tagged-Access Mode (Trusted Classifier) (Continued)

Code Point	Forwarding Class	Loss Priority
ef (010)	best-effort	low
ef1 (011)	fcoe	low
af11 (100)	no-loss	low
af12 (101)	best-effort	low
nc1 (110)	network-control	low
nc2 (111)	network-control	low

Table 49 on page 166 shows the default mapping of IEEE 802.1p code-point values to forwarding classes and loss priorities for ports in access mode (all incoming traffic is mapped to best-effort forwarding classes).

NOTE: Table 49 on page 166 applies only to unicast traffic except on QFX10000 switches and NFX Series devices.

Table 49: Default IEEE 802.1 Classifiers for Ports in Access Mode (Untrusted Classifier)

Code Point	Forwarding Class	Loss Priority
000	best-effort	low
001	best-effort	low
010	best-effort	low
011	best-effort	low

Table 49: Default IEEE 802.1 Classifiers for Ports in Access Mode (Untrusted Classifier) (Continued)

Code Point	Forwarding Class	Loss Priority
100	best-effort	low
101	best-effort	low
110	best-effort	low
111	best-effort	low

Table 50 on page 167 shows the default mapping of IEEE 802.1 code-point values to multdestination (multicast, broadcast, and destination lookup fail traffic) forwarding classes and loss priorities.

NOTE: Table 50 on page 167 does not apply to QFX10000 switches or NFX Series devices.

Table 50: Default IEEE 802.1 Multidestination Classifiers

Code Point	Forwarding Class	Loss Priority
be (000)	mcast	low
be1 (001)	mcast	low
ef (010)	mcast	low
ef1 (011)	mcast	low
af11 (100)	mcast	low
af12 (101)	mcast	low
nc1 (110)	mcast	low

Table 50: Default IEEE 802.1 Multidestination Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
nc2 (111)	mcast	low

Table 51 on page 168 shows the default mapping of DSCP code-point values to forwarding classes and loss priorities for DSCP IP and DCSP IPv6.

NOTE: Table 51 on page 168 applies only to unicast traffic except on QFX10000 switches and NFX Series devices.

Table 51: Default DSCP IP and IPv6 Classifiers

Code Point	Forwarding Class	Loss Priority
ef (101110)	best-effort	low
af11 (001010)	best-effort	low
af12 (001100)	best-effort	low
af13 (001110)	best-effort	low
af21 (010010)	best-effort	low
af22 (010100)	best-effort	low
af23 (010110)	best-effort	low
af31 (011010)	best-effort	low
af32 (011100)	best-effort	low
af33 (011110)	best-effort	low

Table 51: Default DSCP IP and IPv6 Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
af41 (100010)	best-effort	low
af42 (100100)	best-effort	low
af43 (100110)	best-effort	low
be (000000)	best-effort	low
cs1 (001000)	best-effort	low
cs2 (010000)	best-effort	low
cs3 (011000)	best-effort	low
cs4 (100000)	best-effort	low
cs5 (101000)	best-effort	low
nc1 (110000)	network-control	low
nc2 (111000)	network-control	low

NOTE: There are no default DSCP IP or IPv6 multdestination classifiers for multdestination traffic. DSCP IPv6 multdestination classifiers are not supported for multdestination traffic.

Table 52 on page 170 shows the default mapping of MPLS EXP code-point values to forwarding classes and loss priorities, which apply only on QFX10000 switches and NFX Series devices.

Table 52: Default EXP Classifiers on QFX10000 Switches and NFX Series Devices

Code Point	Forwarding Class	Loss Priority
000	best-effort	low
001	best-effort	high
010	expedited-forwarding	low
011	expedited-forwarding	high
100	assured-forwarding	low
101	assured-forwarding	high
110	network-control	low
111	network-control	high

Default Scheduling

The default schedulers allocate egress bandwidth resources to egress traffic as shown in [Table 53 on page 170](#):

Table 53: Default Scheduler Configuration

Default Scheduler and Queue Number	Transmit Rate (Guaranteed Minimum Bandwidth)	Shaping Rate (Maximum Bandwidth)	Excess Bandwidth Sharing	Priority	Buffer Size
best-effort forwarding class scheduler (queue 0)	5% 15% (QFX10000, NFX Series)	None	5% 15% (QFX10000, NFX Series)	low	5% 15% (QFX10000, NFX Series)

Table 53: Default Scheduler Configuration (Continued)

Default Scheduler and Queue Number	Transmit Rate (Guaranteed Minimum Bandwidth)	Shaping Rate (Maximum Bandwidth)	Excess Bandwidth Sharing	Priority	Buffer Size
fcoe forwarding class scheduler (queue 3)	35%	None	35%	low	35%
no-loss forwarding class scheduler (queue 4)	35%	None	35%	low	35%
network-control forwarding class scheduler (queue 7)	5% 15% (QFX10000, NFX Series)	None	5% 15% (QFX10000, NFX Series)	low	5% 15% (QFX10000, NFX Series)
(Excluding QFX10000 and NFX Series) mcast forwarding class scheduler (queue 8)	20%	None	20%	low	20%

NOTE: By default, the minimum guaranteed bandwidth (transmit rate) determines the amount of excess (extra) bandwidth that a queue can share. Extra bandwidth is allocated to queues in proportion to the transmit rate of each queue. On switches that support the `excess-rate` statement, you can override the default setting and configure the excess bandwidth percentage independently of the transmit rate on queues that are not strict-high priority queues.

By default, only the four (QFX10000 switches and NFX Series devices) or five (other switches) default schedulers shown in [Table 53 on page 170](#) have traffic mapped to them. Only the forwarding classes and queues associated with the default schedulers receive default bandwidth, based on the default scheduler transmit rate. (You can configure schedulers and forwarding classes to allocate bandwidth to other queues or to change the bandwidth and other scheduling properties of a default queue.)

On QFX10000 switches and NFX Series devices, if a forwarding class does not transport traffic, the bandwidth allocated to that forwarding class is available to other forwarding classes. Unicast and multdestination (multicast, broadcast, and destination lookup fail) traffic use the same forwarding classes and output queues.

On switches other than QFX10000 and NFX Series devices, multidestination queue 11 receives enough bandwidth from the default multidestination scheduler to handle CPU-generated multidestination traffic.

On QFX10000 and NFX Series devices, default scheduling is port scheduling. Default hierarchical scheduling, known as enhanced transmission selection (ETS, defined in IEEE 802.1Qaz), allocates the total port bandwidth to the four default forwarding classes served by the four default schedulers, as defined by the four default schedulers. The result is the same as direct port scheduling. Configuring hierarchical port scheduling, however, enables you to group forwarding classes that carry similar types of traffic into forwarding class sets (also called priority groups), and to assign port bandwidth to each forwarding class set. The port bandwidth assigned to the forwarding class set is then assigned to the forwarding classes within the forwarding class set. This hierarchy enables you to control port bandwidth allocation with greater granularity, and enables hierarchical sharing of extra bandwidth to better utilize link bandwidth.

Except on QFX10000 switches and NFX Series devices, default hierarchical scheduling divides the total port bandwidth between two groups of traffic: unicast traffic and multidestination traffic. By default, unicast traffic consists of queue 0 (best-effort forwarding class), queue 3 (fcfe forwarding class), queue 4 (no-loss forwarding class), and queue 7 (network-control forwarding class). Unicast traffic receives and shares a total of 80 percent of the port bandwidth. By default, multidestination traffic (mcast queue 8) receives a total of 20 percent of the port bandwidth. So on a 10-Gigabit port, unicast traffic receives 8-Gbps of bandwidth and multidestination traffic receives 2-Gbps of bandwidth.

NOTE: Except on QFX5200, QFX5210, and QFX10000 switches and NFX Series devices, which do not support queue 11, multidestination queue 11 also receives a small amount of default bandwidth from the multidestination scheduler. CPU-generated multidestination traffic uses queue 11, so you might see a small number of packets egress from queue 11. In addition, in the unlikely case that firewall filter match conditions map multidestination traffic to a unicast forwarding class, that traffic uses queue 11.

Default scheduling uses weighted round-robin (WRR) scheduling. Each queue receives a portion (weight) of the total available interface bandwidth. The scheduling weight is based on the transmit rate of the default scheduler for that queue. For example, queue 7 receives a default scheduling weight of 5 percent, or 15 percent on QFX10000 and NFX Series devices, of the available bandwidth, and queue 4 receives a default scheduling weight of 35 percent of the available bandwidth. Queues are mapped to forwarding classes, so forwarding classes receive the default bandwidth for the queues to which they are mapped.

On QFX10000 switches and NFX Series devices, for example, queue 7 is mapped to the network-control forwarding class and queue 4 is mapped to the no-loss forwarding class. Each forwarding class receives the default bandwidth for the queue to which it is mapped. Unused bandwidth is shared with other default queues.

If you want non-default (unconfigured) queues to forward traffic, you should explicitly map traffic to those queues (configure the forwarding classes and queue mapping) and create schedulers to allocate bandwidth to those queues. By default, queues 1, 2, 5, and 6 are unconfigured.

Except on QFX5200, QFX5210, and QFX10000 switches and NFX Series devices, which do not support them, multidestination queues 9, 10, and 11 are unconfigured. Unconfigured queues have a default scheduling weight of 1 so that they can receive a small amount of bandwidth in case they need to forward traffic. However, queue 11 can use more of the default multidestination scheduler bandwidth if necessary to handle CPU-generated multidestination traffic.

NOTE: All four (two on QFX5200 and QFX5210 switches) multidestination queues have a scheduling weight of 1. Because by default multidestination traffic goes to queue 8, queue 8 receives almost all of the multidestination bandwidth. (There is no traffic on queue 9 and queue 10, and very little traffic on queue 11, so there is almost no competition for multidestination bandwidth.)

However, if you explicitly configure queue 9, 10, or 11 (by mapping code points to the unconfigured multidestination forwarding classes using the multidestination classifier), the explicitly configured queues share the multidestination scheduler bandwidth equally with default queue 8, because all of the queues have the same scheduling weight (1). To ensure that multidestination bandwidth is allocated to each queue properly and that the bandwidth allocation to the default queue (8) is not reduced too much, we strongly recommend that you configure a scheduler if you explicitly classify traffic into queue 9, 10, or 11.

If you map traffic to an unconfigured queue, the queue receives only the amount of excess bandwidth proportional to its default weight (1). The actual amount of bandwidth an unconfigured queue gets depends on how much bandwidth the other queues are using.

If some queues use less than their allocated amount of bandwidth, the unconfigured queues can share the unused bandwidth. Sharing unused bandwidth is one of the key advantages of hierarchical port scheduling. Configured queues have higher priority for bandwidth than unconfigured queues, so if a configured queue needs more bandwidth, then less bandwidth is available for unconfigured queues. Unconfigured queues always receive a minimum amount of bandwidth based on their scheduling weight (1). If you map traffic to an unconfigured queue, to allocate bandwidth to that queue, configure a scheduler for the forwarding class that is mapped to the queue.

Default DCBX Advertisement

When you configure hierarchical scheduling on an interface, DCBX advertises each priority group, the priorities in each priority group, and the bandwidth properties of each priority and priority group.

If you do not configure hierarchical scheduling on an interface, DCBX advertises the automatically created default priority group and its priorities. DCBX also advertises the default bandwidth allocation of the priority group, which is 100 percent of the port bandwidth.

Default Scheduling and Classification Summary

If you do not configure scheduling on an interface:

- Default classifiers classify ingress traffic.
- Default schedulers schedule egress traffic.
- DCBX advertises a single default priority group with 100 percent of the port bandwidth allocated to that priority group. All priorities (forwarding classes) are assigned to the default priority group and receive bandwidth based on their default schedulers. The default priority group is generated automatically and is not user-configurable.

RELATED DOCUMENTATION

Understanding CoS Packet Flow

Understanding CoS Hierarchical Port Scheduling (ETS)

Understanding Default CoS Settings

Understanding CoS Virtual Output Queues (VOQs) on QFX10000 Switches

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

Understanding DCB Features and Requirements

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

Example: Configuring Unicast Classifiers

Example: Configuring Queue Schedulers

Understanding CoS Scheduling on QFabric System Node Device Fabric (fte) Ports

IN THIS SECTION

- [Hierarchical Scheduling Architecture on QFabric System Node Devices | 175](#)

- [Default Scheduling on Node Device Fabric Interfaces | 176](#)
- [Configuring Scheduling on Node Device Fabric Interfaces | 178](#)

Beginning with Junos OS Release 13.1, you can configure two-tier hierarchical scheduling (enhanced transmission selection, IEEE 802.1Qaz) on the fabric (fte) ports of QFabric system Node devices. Configuring CoS on Node device fabric interfaces provides increased control over traffic scheduling and helps to ensure predictable bandwidth consumption.

You can configure CoS on the following QFabric system interface types:

- Node device access interfaces (xe interfaces)—Schedule traffic on the output queues of the 10-Gigabit Ethernet access ports using standard Node device CoS scheduling configuration components, as described elsewhere in the QFX Series documentation. You can configure different scheduling for different ports and output queues.
- Node device fabric interfaces (fte interfaces)—Schedule traffic on the output queues of the 40-Gbps fabric interfaces that connect a Node device to a QFX3008-I or a QFX3600-I Interconnect device using standard Node device CoS scheduling configuration components. You can configure different scheduling for different interfaces and output queues.

This topic describes:

Hierarchical Scheduling Architecture on QFabric System Node Devices

CoS architecture on Node device access interfaces is the same as CoS architecture on standalone switch access interfaces. CoS architecture on Node device fabric interfaces is also the same as the CoS architecture on the access interfaces. You apply schedulers to queues (priorities), fc-sets (priority groups), and interfaces in the same hierarchical manner as described in [Understanding CoS Hierarchical Port Scheduling \(ETS\)](#).

You configure scheduling on Node device fabric interfaces (fte interfaces) using the same statements and configuration constructs that you use to configure scheduling on Node device access interfaces (xe interfaces). For example, on Node device fabric interfaces you can:

- Define up to four fc-sets (three unicast, one multidestination)

NOTE: If the fabric interface handles strict-high priority traffic, you must define a separate fc-set (priority group) for strict-high priority traffic. Strict-high priority traffic cannot be mixed

with traffic of other priorities in an fc-set. For example, you might choose to create different fc-sets for best-effort, lossless, strict-high priority, and multidestination traffic.

- Map forwarding classes to fc-sets
- Configure scheduling for each forwarding class (scheduler)
- Configure scheduling for each fc-set (traffic control profile)

The differences in configuring CoS on Node device fabric interfaces compared to configuring CoS on Node device access interfaces are:

- You specify a Node device *fabric* interface instead of a Node device *access* interface when you apply CoS to an interface.
- You cannot attach classifiers, congestion notification profiles, or rewrite rules to fabric interfaces. Also, you cannot configure buffer settings on fabric interfaces. You can only attach fc-sets and traffic control profiles.

Default Scheduling on Node Device Fabric Interfaces

Default scheduling on Node device fabric interfaces is the same as default scheduling on Node device access interfaces. Only the default forwarding classes (best-effort, network-control, fcoe, no-loss, and multidestination) receive port bandwidth, based on the default minimum guaranteed bandwidth (transmit rate) scheduler settings for each default forwarding class.

To transport traffic on Node device fabric interfaces, the system organizes the default forwarding classes into three *class groups*. Class groups are not user-configurable. The three class groups are:

- **Unicast**—All traffic in the default forwarding classes best-effort, network-control, fcoe, and no-loss belong to this default class group.
- **Multidestination**—All traffic in the default forwarding class mcast belongs to this default class group.
- **Strict-high priority**—There is no default strict-high priority forwarding class, so there is no default strict-high priority class group and there is no default configuration for strict-high priority traffic.

NOTE: If you configure strict-high priority forwarding classes, you must also configure an fc-set (priority group) for strict-high priority traffic, map the strict-high priority forwarding classes to the strict-high priority fc-set, create a scheduler for the strict-high priority traffic and map it to the strict-high priority forwarding classes, create a traffic control profile for the strict-high priority traffic, and apply the strict-high priority fc-set and traffic control profile to the appropriate fabric interfaces.

The default forwarding classes receive port bandwidth based on their default transmit rate settings (weights). Forwarding classes that are not default forwarding classes receive no default bandwidth.

Default class group scheduling uses weighted round-robin (WRR) scheduling, in which each class group receives a portion of the total available fabric interface bandwidth based on the class group traffic type, as shown in [Table 54 on page 177](#). Within each class group, the scheduler bandwidth allocation for individual forwarding classes is based on the default transmit rate for each forwarding class.

Table 54: Class Group Default Scheduling Properties and Membership on Node Device Fabric Interfaces

Class Group	Forwarding Class Mapping and Bandwidth Allocation (Default Transmit Rate)	Class Group Scheduling Properties (Weight)
Unicast	<ul style="list-style-type: none"> • best-effort (5%) • fcoe (35%) • no-loss (35%) • network-control (5%) 	Traffic in the unicast class group receives an 80% weight in the weighted round-robin (WRR) calculations. After the strict-high priority class group has been served, the unicast class group receives 80% of the remaining fabric bandwidth. (If more bandwidth is available, the unicast class group can use more bandwidth.)
Multidestination	<ul style="list-style-type: none"> • mcast (20%) 	Traffic in the multidestination class group receives a 20% weight in the WRR calculations. After the strict-high priority class group has been served, the multidestination class group receives 20% of the remaining fabric bandwidth. (If more bandwidth is available, the multidestination class group can use more bandwidth.)

NOTE: Strict-high priority traffic is served first, before any other traffic is served. Strict-high priority traffic receives all of the bandwidth it needs to empty its queues and therefore can starve other types of traffic during periods of high-volume strict-high priority traffic. Plan carefully and use caution when determining how much traffic to configure as strict-high priority traffic. We recommend that you always configure a shaping rate in the strict-high priority scheduler to set a maximum bandwidth limit for strict-high priority traffic.

Configuring Scheduling on Node Device Fabric Interfaces

If you do not want to use default CoS scheduling on Node device fabric interfaces, you can configure two-tier hierarchical scheduling (ETS) the same way that you configure ETS on Node device access interfaces.

Similarities Between Node Device Fabric Interface and Access Interface Scheduling

Configuring scheduling on a Node device fabric interface is similar to configuring scheduling on an access interface in many ways. In both cases, you configure:

- Schedulers to specify the output scheduling for forwarding class traffic
- Scheduler maps to map schedulers to forwarding classes
- Forwarding classes (or use the default forwarding classes)
- Forwarding class sets (groups of forwarding classes that require similar CoS treatment)
- A separate fc-set for strict-high priority traffic (an fc-set cannot contain a mix of strict-high priority traffic and traffic with a different priority)
- Traffic control profiles to specify the output scheduling for fc-sets
- Traffic control profile and fc-set mapping to interfaces

On Node device fabric interfaces, you configure ETS in the same way, and ETS works the same way as on Node device access interfaces

In addition, strict-high priority queues are served first, and then the remaining port bandwidth is allocated to other traffic. Unless you configure a shaping rate in the scheduler for strict-high priority traffic, a strict-high priority queue can consume all of the port bandwidth and starve other queues, so we recommend that you always configure a shaping rate on strict-high priority traffic.

Differences Between Node Device Fabric Interface and Access Interface Scheduling

Configuring scheduling on a Node device fabric interface differs from configuring scheduling on an access interface in several ways. On fabric interfaces:

- You cannot attach classifiers.
- You cannot attach congestion notification profiles (flow control is applied automatically to lossless forwarding classes).
- You cannot attach rewrite rules.
- You cannot configure buffer settings.

- You specify a Node device fabric interface name instead of a Node device access interface name when you apply CoS to an interface.

RELATED DOCUMENTATION

Understanding CoS Fabric Forwarding Class Sets

Understanding CoS Output Queue Schedulers

Understanding CoS Hierarchical Port Scheduling (ETS)

Understanding Default CoS Settings

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

[Example: Configuring CoS Scheduling Across the QFabric System | 353](#)

Example: Configuring Queue Schedulers

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring WRED Drop Profiles

Example: Configuring Drop Profile Maps

Understanding CoS Scheduling Behavior and Configuration Considerations

Many factors affect scheduling configuration and bandwidth requirements, including:

- When you configure bandwidth for a forwarding class (each forwarding class is mapped to a queue) or a forwarding class set (priority group), the switch considers only the data as the configured bandwidth. The switch does not account for the bandwidth consumed by the preamble and the interframe gap (IFG). Therefore, when you calculate and configure the bandwidth requirements for a forwarding class or for a forwarding class set, consider the preamble and the IFG as well as the data in the calculations.
- When you configure a forwarding class to carry traffic on the switch (instead of using only default forwarding classes), you must also define a scheduling policy for the user-configured forwarding class. Some switches support enhanced transmission selection (ETS) hierarchical port scheduling, some switches support direct port scheduling, and some switches support both methods of scheduling.

For ETS hierarchical port scheduling, defining a hierarchical scheduling policy using ETS means:

- Mapping a scheduler to the forwarding class in a scheduler map

- Including the forwarding class in a forwarding class set
- Associating the scheduler map with a traffic control profile
- Attaching the traffic control profile to a forwarding class set and an interface

On switches that support port scheduling, defining a scheduling policy means:

- Mapping a scheduler to the forwarding class in a scheduler map.
- Applying the scheduler map to one or more interfaces.
- On each physical interface, either all forwarding classes that are being used on the interface must have rewrite rules configured, or no forwarding classes that are being used on the interface can have rewrite rules configured. On any physical port, do not mix forwarding classes with rewrite rules and forwarding classes without rewrite rules.
- For packets that carry both an inner VLAN tag and an outer VLAN tag, rewrite rules rewrite only the outer VLAN tag.
- For ETS hierarchical port scheduling, configuring the minimum guaranteed bandwidth (`transmit-rate`) for a forwarding class does not work unless you also configure the minimum guaranteed bandwidth (`guaranteed-rate`) for the forwarding class set in the traffic control profile.

Additionally, the sum of the transmit rates of the forwarding classes in a forwarding class set should not exceed the guaranteed rate for the forwarding class set. (You cannot guarantee a minimum bandwidth for the queues that is greater than the minimum bandwidth guaranteed for the entire set of queues.) If you configure transmit rates whose sum exceeds the guaranteed rate of the forwarding class set, the commit check fails and the system rejects the configuration.

- For ETS hierarchical port scheduling, the sum of the forwarding class set guaranteed rates cannot exceed the total port bandwidth. If you configure guaranteed rates whose sum exceeds the port bandwidth, the system sends a syslog message to notify you that the configuration is not valid. However, the system does not perform a commit check. If you commit a configuration in which the sum of the guaranteed rates exceeds the port bandwidth, the hierarchical scheduler behaves unpredictably.
- For ETS hierarchical port scheduling, if you configure the `guaranteed-rate` of a forwarding class set as a percentage, configure all of the transmit rates associated with that forwarding class set as percentages. In this case, if any of the transmit rates are configured as absolute values instead of percentages, the configuration is not valid and the system sends a syslog message.
- There are several factors to consider if you want to configure a strict-high priority queue (forwarding class):
 - On QFX5200, QFX3500, and QFX3600 switches and on QFabric systems, you can configure only one strict-high priority queue (forwarding class).

On QFX5100 and EX4600 switches, you can configure only one forwarding-class-set (priority group) as strict-high priority. All queues which are part of that strict-high forwarding class set then act as strict-high queues.

On QFX10000 switches, there is no limit to the number of strict-high priority queues you can configure.

- You cannot configure a minimum guaranteed bandwidth (`transmit-rate`) for a strict-high priority queue on QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems.

On QFX5200 and QFX10000 switches, you can set the `transmit-rate` on strict-high priority queues to set a limit on the amount of traffic that the queue treats as strict-high priority traffic. Traffic in excess of the `transmit-rate` is treated as best-effort traffic, and receives an excess bandwidth sharing weight of “1”, which is the proportion of extra bandwidth the strict-high priority queue can share on the port. Queues that are not strict-high priority queues use the `transmit rate` (default) or the configured excess rate to determine the proportion (weight) of extra port bandwidth the queue can share. However, you cannot configure an excess rate on a strict-high priority queue, and you cannot change the excess bandwidth sharing weight of “1” on a strict-high priority queue.

For ETS hierarchical port scheduling, you cannot configure a minimum guaranteed bandwidth (`guaranteed-rate`) for a forwarding class set that includes a strict-high priority queue.

- Except on QFX10000 switches, for ETS hierarchical port scheduling only, you must create a separate forwarding class set for a strict-high priority queue. On QFX10000 switches, you can mix strict-high priority and low priority queues in the same forwarding class set.
- Except on QFX10000 switches, for ETS hierarchical port scheduling, only one forwarding class set can contain a strict-high priority queue. On QFX10000 switches, this restriction does not apply.
- Except on QFX10000 switches, for ETS hierarchical port scheduling, a strict-high priority queue cannot belong to the same forwarding class set as queues that are not strict-high priority. (You cannot mix a strict-high priority forwarding class with forwarding classes that are not strict-high priority in one forwarding class set.) On QFX10000 switches, you can mix strict-high priority and low priority queues in the same forwarding class set.
- For ETS hierarchical port scheduling on switches that use different forwarding class sets for unicast and multdestination (multicast, broadcast, and destination lookup fail) traffic, a strict-high priority queue cannot belong to a multdestination forwarding class set.
- On QFX10000 systems, we recommend that you always configure a transmit rate on strict-high priority queues to prevent them from starving other queues. If you do not apply a transmit rate to limit the amount of bandwidth strict-high priority queues can use, then strict-high priority queues can use all of the available port bandwidth and starve other queues on the port.

On QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, we recommend that you always apply a shaping rate to the strict-high priority queue to prevent it from starving other queues. If you do not apply a shaping rate to limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

- On QFabric systems, if any queue that contains outgoing packets does not transmit packets for 12 consecutive seconds, the port automatically resets. Failure of a queue to transmit packets for 12 consecutive seconds might be due to:
 - A strict-high priority queue consuming all of the port bandwidth
 - Several queues consuming all of the port bandwidth
 - Any queue or port receiving continuous *priority-based flow control* (PFC) or 802.3x Ethernet PAUSE messages (received PFC and PAUSE messages prevent a queue or a port, respectively, from transmitting packets because of network congestion)
 - Other conditions that prevent a queue from obtaining port bandwidth for 12 consecutive seconds

If the cause is a strict-high priority queue consuming all of the port bandwidth, use rate shaping to configure a maximum rate for the strict-high priority queue and prevent it from using all of the port bandwidth. To configure rate shaping, include the `shaping-rate (rate | percent percentage)` statement at the `[edit class-of-service schedulers scheduler-name]` hierarchy level and apply the shaping rate to the strict-high priority scheduler. We recommend that you always apply a shaping rate to strict-high priority traffic to prevent the strict-high priority queue from starving other queues.

If several queues consume all of the port bandwidth, you can use a scheduler to rate shape those queues and prevent them from using all of the port bandwidth.

- For transmit rates below 1 Gbps, we recommend that you configure the transmit rate as a percentage instead of as a fixed rate. This is because the system converts fixed rates into percentages and might round small fixed rates to a lower percentage. For example, a fixed rate of 350 Mbps is rounded down to 3 percent instead of 3.5 percent.
- When you set the maximum bandwidth for a queue or for a priority group (`shaping-rate`) at 100 Kbps or lower, the traffic shaping behavior is accurate only within ± 20 percent of the configured shaping rate.
- On QFX10000 switches, configuring rate shaping (`[set class-of-service schedulers scheduler-name transmit-rate (rate | percentage) exact]`) on a LAG interface using the `[edit class-of-service interfaces lag-interface-name scheduler-map scheduler-map-name]` statement can result in scheduled traffic streams receiving more LAG link bandwidth than expected.

You configure rate shaping in a scheduler to set the maximum bandwidth for traffic assigned to a forwarding class on a particular output queue on a port. For example, you can use a scheduler to

configure rate shaping on traffic assigned to the best-effort forwarding class mapped to queue 0, and then apply the scheduler to an interface using a scheduler map, to set the maximum bandwidth for best-effort traffic mapped to queue 0 on that port. Traffic in the best-effort forwarding can use no more than the amount of port bandwidth specified by the transmit rate when you use the `exact` option.

LAG interfaces are composed of two or more Ethernet links bundled together to function as a single interface. The switch can hash traffic entering a LAG interface onto any member link in the LAG interface. When you configure rate shaping and apply it to a LAG interface, the way that the switch applies the rate shaping to traffic depends on how the switch hashes the traffic onto the LAG links.

To illustrate how link hashing affects the way the switch applies a shaping rate to LAG traffic, let's look at a LAG interface (`ae0`) that has two member links (`xe-0/0/20` and `xe-0/0/21`). On LAG `ae0`, we configure rate shaping of 2g for traffic assigned to the best-effort forwarding class, which is mapped to output queue 0. When traffic in the best-effort forwarding class reaches the LAG interface, the switch hashes the traffic onto one of the two member links.

If the switch hashes all of the best-effort traffic onto the same LAG link, the traffic receives a maximum of 2g bandwidth on that link. In this case, the intended cumulative limit of 2g for best-effort traffic on the LAG is enforced.

However, if the switch hashes the best-effort traffic onto both of the LAG links, the traffic receives a maximum of 2g bandwidth on *each* LAG link, not 2g as a cumulative total for the entire LAG, so the best-effort traffic receives a maximum of 4g on the LAG, not the 2g set by the rate shaping configuration. When hashing spreads the traffic assigned to an output queue (which is mapped to a forwarding class) across multiple LAG links, the effective rate shaping (cumulative maximum bandwidth) on the LAG is:

(number of LAG member interfaces) x (rate shaping for the output queue) = cumulative LAG rate shaping

- On switches that do not use virtual output queues (VOQs), ingress port congestion can occur during periods of egress port congestion if an ingress port forwards traffic to more than one egress port, and at least one of those egress ports experiences congestion. If this occurs, the congested egress port can cause the ingress port to exceed its fair allocation of ingress buffer resources. When the ingress port exceeds its buffer resource allocation, frames are dropped at the ingress. Ingress port frame drop affects not only the congested egress ports, but also all of the egress ports to which the congested ingress port forwards traffic.

If a congested ingress port drops traffic that is destined for one or more uncongested egress ports, configure a weighted random early detection (WRED) drop profile and apply it to the egress queue that is causing the congestion. The drop profile prevents the congested egress queue from affecting egress queues on other ports by dropping frames at the egress instead of causing congestion at the ingress port.

NOTE: On systems that support lossless transport, do not configure drop profiles for lossless forwarding classes such as the default `fcoe` and `no-loss` forwarding classes. FCoE and other lossless traffic queues require lossless behavior. Use priority-based flow control (PFC) to prevent frame drop on lossless priorities.

- On systems that use different classifiers for unicast and multdestination traffic and that support lossless transport, on an ingress port, do not configure classifiers that map the same IEEE 802.1p code point to both a multdestination traffic flow and a lossless unicast traffic flow (such as the default lossless `fcoe` or `no-loss` forwarding classes). Any code point used for multdestination traffic on a port should not be used to classify unicast traffic into a lossless forwarding class on the same port.

If a multdestination traffic flow and a lossless unicast traffic flow use the same code point on a port, the multdestination traffic is treated the same way as the lossless traffic. For example, if priority-based flow control (PFC) is applied to the lossless traffic, the multdestination traffic of the same code point is also paused. During periods of congestion, treating multdestination traffic the same as lossless unicast traffic can create ingress port congestion for the multdestination traffic and affect the multdestination traffic on all of the egress ports the multdestination traffic uses.

For example, the following configuration can cause ingress port congestion for the multdestination flow:

- For unicast traffic, IEEE 802.1p code point 011 is classified into the `fcoe` forwarding class:

```
user@switch# set class-of-service classifiers ieee-802.1 ucast_cl forwarding-class fcoe
loss-priority low code-points 011
```

- For multdestination traffic, IEEE 802.1p code point 011 is classified into the `mcast` forwarding class:

```
user@switch# set class-of-service classifiers ieee-802.1 mcast-cl forwarding-class mcast
loss-priority low code-points 011
```

- The unicast classifier that maps traffic with code point 011 to the `fcoe` forwarding class is mapped to interface `xe-0/0/1`:

```
user@switch# set class-of-service interfaces xe-0/0/1 unit 0 classifiers ieee-802.1
ucast_cl
```

4. The multdestination classifier that maps traffic with code point 011 to the mcast forwarding class is mapped to all interfaces (multidestination traffic maps to all interfaces and cannot be mapped to individual interfaces):

```
user@switch# set class-of-service multi-destination classifiers ieee-802.1 mcast-cl
```

Because the same code point (011) maps unicast traffic to a lossless traffic flow and also maps multidestination traffic to a multidestination traffic flow, the multidestination traffic flow might experience ingress port congestion during periods of congestion.

To avoid ingress port congestion, do not map the code point used by the multidestination traffic to lossless unicast traffic. For example:

1. Instead of classifying code point 011 into the fcoe forwarding class, classify code point 011 into the best-effort forwarding class:

```
user@switch# set class-of-service classifiers ieee-802.1 ucast_cl forwarding-class best-effort loss-priority low code-points 011
```

2.

```
user@switch# set class-of-service classifiers ieee-802.1 mcast-cl forwarding-class mcast loss-priority low code-points 011
```

3.

```
user@switch# set class-of-service interfaces xe-0/0/1 unit 0 classifiers ieee-802.1 ucast_cl
```

4.

```
user@switch# set class-of-service multi-destination classifiers ieee-802.1 mcast-cl
```

Because the code point 011 does not map unicast traffic to a lossless traffic flow, the multidestination traffic flow does not experience ingress port congestion during periods of congestion.

The best practice is to classify unicast traffic with IEEE 802.1p code points that are also used for multidestination traffic into best-effort forwarding classes.

Understanding CoS Output Queue Schedulers

IN THIS SECTION

- [Output Queue Scheduling Components | 187](#)
- [Default Schedulers | 189](#)
- [Transmit Rate \(Minimum Guaranteed Bandwidth\) | 189](#)
- [Sharing Extra Bandwidth | 190](#)
- [Shaping Rate \(Maximum Bandwidth\) | 191](#)
- [Scheduling Priority | 191](#)
- [Scheduler Drop-Profile Maps | 191](#)
- [Buffer Size | 192](#)
- [Explicit Congestion Notification | 193](#)
- [Scheduler Maps | 193](#)

Output queue scheduling defines the class-of-service (CoS) properties of output queues. Output queues are mapped to forwarding classes, and classifiers map incoming traffic into forwarding classes based on IEEE 802.1p or DSCP code points. Output queue properties include the amount of interface bandwidth assigned to the queue, the size of the memory buffer allocated for storing packets, the priority of the queue, and the weighted random early detection (WRED) drop profiles associated with the queue. Queue scheduling works with priority group scheduling to create a two-tier hierarchical scheduler.

The hierarchical scheduler allocates port bandwidth to a group of queues (forwarding classes) called a priority group (forwarding class set), and queue scheduling determines the portion of the priority group's bandwidth that a particular queue can use. So the first scheduling tier is allocating port bandwidth to a forwarding class set, and the second scheduling tier is allocating forwarding class set bandwidth to forwarding classes (queues).

Scheduler maps associate queue schedulers with forwarding classes. The queue mapped to a forwarding class receives the scheduling resources assigned to that forwarding class. You associate a scheduler map with a traffic control profile, and then associate the traffic control profile with a forwarding class set (priority group) and a port interface to apply scheduling to a port. In conjunction with the priority group scheduling configured in the traffic control profile, queue scheduling configures the packet schedulers and weighted random early detection (WRED) packet drop processes for queues.

NOTE: When you configure bandwidth for a queue or a priority group, the switch considers only the data as the configured bandwidth. The switch does not account for the bandwidth consumed by the preamble and the interframe gap (IFG). Therefore, when you calculate and configure the bandwidth requirements for a queue or for a priority group, consider the preamble and the IFG as well as the data in the calculations.

Output Queue Scheduling Components

[Table 55 on page 187](#) provides a quick reference to the scheduler components you can configure to determine the bandwidth properties of output queues (forwarding classes), and [Table 56 on page 188](#) provides a quick reference to some related scheduling configuration components.

Table 55: Output Queue Scheduler Components

Output Queue Scheduler Component	Description
Buffer size	<p>Sets the size of the queue buffer.</p> <p>See Understanding CoS Buffer Configuration.</p>
Drop profile map	<p>Maps a drop profile to a loss priority. Drop profile map components include:</p> <ul style="list-style-type: none"> Drop profile—Sets the probability of dropping packets as the queue fills up. Loss priority—Sets the traffic loss priority to which a drop profile applies. <p>See Configuring CoS Drop Profile Maps.</p>
Explicit congestion notification	<p>Enables explicit congestion notification (ECN) on the queue.</p> <p>See Understanding CoS Explicit Congestion Notification.</p>
Priority	<p>Sets the scheduling priority applied to the queue.</p> <p>See Defining CoS Queue Scheduling Priority.</p>

Table 55: Output Queue Scheduler Components *(Continued)*

Output Queue Scheduler Component	Description
Shaping rate	<p>Sets the maximum bandwidth the queue can consume.</p> <p>TIP: On QFX5200 Series switches, a granularity of 64kbps is supported for the shaping rate.</p> <p>See Understanding CoS Priority Group Shaping and Queue Shaping (Maximum Bandwidth).</p>
Transmit rate	<p>Sets the minimum guaranteed bandwidth for the queue. Extra bandwidth is shared among queues in proportion to the minimum guaranteed bandwidth of each queue. See Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth.</p>

Table 56: Other Scheduling Components

Other Scheduling Components	Description
Forwarding class	<p>Maps traffic to an output queue. Classifiers map forwarding classes to IEEE 802.1p, DSCP, or EXP code points. A forwarding class, an output queue, and code point bits are mapped to each other and identify the same traffic. (The code point bits identify incoming traffic. Classifiers assign traffic to forwarding classes based on the code point bits. Forwarding classes are mapped to output queues. This mapping determines the output queue each class of traffic uses on the switch egress interfaces.)</p>
Output queue	<p>Buffers traffic before the switch forwards the traffic out the egress interface. Output queues are mapped to forwarding classes. The switch applies CoS properties defined in schedulers to output queues, by mapping forwarding classes to schedulers in scheduler maps. The queue mapped to the forwarding class has the CoS properties defined in the scheduler mapped to that forwarding class.</p>

Table 56: Other Scheduling Components (Continued)

Other Scheduling Components	Description
Scheduler map	Maps schedulers to forwarding classes (forwarding classes are mapped to queues, so a forwarding class represents a queue, and the scheduler mapped to a forwarding class determines the CoS properties of the output queue mapped to that forwarding class).
Traffic control profile	Configures scheduling for the forwarding class set (priority group), and associates a scheduler map with the forwarding class set to apply queue scheduling to the forwarding classes in the forwarding class set. Extra port bandwidth is shared among forwarding class sets in proportion to the minimum guaranteed bandwidth of each forwarding class set.
Forwarding class set	Name of a priority group. You map forwarding classes to forwarding class sets. A forwarding class set consists of one or more forwarding classes.

Default Schedulers

Each forwarding class requires a scheduler to set the CoS properties of the forwarding class and its output queue. You can use the default schedulers or you can define new schedulers for the associated forwarding classes. For any other forwarding class, you must explicitly configure a scheduler. For more information, see [Default Scheduling](#).

Transmit Rate (Minimum Guaranteed Bandwidth)

The transmit rate determines the minimum guaranteed bandwidth for each forwarding class. The switch applies the minimum bandwidth guarantee to the output queue mapped to the forwarding class. The transmit rate also determines how much excess (extra) bandwidth each low-priority queue can share; each queue shares extra bandwidth in proportion to its transmit rate. You specify the rate in bits per second as a fixed value such as 1 Mbps or as a percentage of the total forwarding class set minimum guaranteed bandwidth (the guaranteed rate set in the traffic control profile). Either the default scheduler or a scheduler you configure allocates a portion of the outgoing interface bandwidth to each forwarding class in proportion to the transmit rate.

NOTE: For transmit rates below 1 Gbps, we recommend that you configure the transmit rate as a percentage instead of as a fixed rate. This is because the system converts fixed rates into percentages and may round small fixed rates to a lower percentage. For example, a fixed rate of 350 Mbps is rounded down to 3 percent.

You cannot configure a transmit rate for a strict-high priority queue. Queues with a configured transmit rate cannot be included in a forwarding class set that has a strict-high priority queue (you cannot mix strict-high priority queues and queues that are not strict-high priority in the same forwarding class set).

The allocated bandwidth can exceed the configured minimum rate if additional bandwidth is available from other queues in the forwarding class set that are not using all of their allocated bandwidth. During periods of congestion, the configured transmit rate is the guaranteed bandwidth minimum for the queue. This behavior enables you to ensure that each queue receives the amount of bandwidth appropriate to its level of service and is also able to share unused bandwidth.

NOTE: Configuring the minimum guaranteed bandwidth (transmit rate) for a forwarding class does not work unless you also configure the minimum guaranteed bandwidth (guaranteed rate) for the forwarding class set in the traffic control profile.

Additionally, the sum of the transmit rates of the queues in a forwarding class set should not exceed the guaranteed rate for the forwarding class set. (You cannot guarantee a combined minimum bandwidth for the queues that is greater than the minimum bandwidth guaranteed for the entire set of queues.)

For more information, see [Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth](#).

Sharing Extra Bandwidth

Extra bandwidth is available to low-priority queues when a forwarding class set does not use its full amount of minimum guaranteed bandwidth (guaranteed-rate). Extra bandwidth is shared among the forwarding classes in a forwarding class set in proportion to the minimum guaranteed bandwidth (transmit-rate) of each queue.

For example, in a forwarding class set, Queue A has a transmit rate of 1 Gbps, Queue B has a transmit rate of 1 Gbps, and Queue C has a transmit rate of 2 Gbps. After servicing the minimum guaranteed bandwidth of these queues, the forwarding class set has an extra 2 Gbps of bandwidth available, and all three queues still have packets to forward. The queues receive the extra bandwidth in proportion to their transmit rates, so Queue A receives an extra 500 Mbps, Queue B receives an extra 500 Mbps, and Queue C receives an extra 1 Gbps.

Shaping Rate (Maximum Bandwidth)

The shaping rate sets the maximum bandwidth that a forwarding class can consume. You specify the rate in bits per second as a fixed value, such as 3 Mbps or as a percentage of the total forwarding class set maximum bandwidth (the shaping rate set in the traffic control profile).

The maximum bandwidth for a queue depends on the total bandwidth available to the forwarding class set to which the queue belongs, and on how much bandwidth the other queues in the forwarding class set consume.

NOTE: On QFabric systems, if any queue that contains outgoing packets does not transmit packets for 12 consecutive seconds, the port automatically resets. A strict-high priority queue (or several queues with higher priorities than the starved queue) can consume all of the port bandwidth and prevent another queue from transmitting packets. To prevent a queue from being starved for bandwidth, you can configure a shaping rate on the queue or queues to prevent them from consuming all of the port bandwidth.

NOTE: We recommend that you always configure a shaping rate in the scheduler for strict-high priority queues to prevent them from starving other queues.

For more information, see [Understanding CoS Priority Group Shaping and Queue Shaping \(Maximum Bandwidth\)](#).

Scheduling Priority

Scheduling priority determines the order in which an interface transmits traffic from its output queues. This ensures that queues containing important traffic receive prioritized access to the outgoing interface bandwidth. The priority setting in the scheduler determines the priority for the queue.

For more information, see [Defining CoS Queue Scheduling Priority](#).

Scheduler Drop-Profile Maps

Drop-profile maps associate drop profiles with queue schedulers and packet loss priorities (PLPs). Drop profiles set thresholds for dropping packets during periods of congestion, based on the queue fill level and a percentage probability of dropping packets at the specified queue fill level. At different fill levels, a drop profile sets different probabilities of dropping a packet during periods of congestion.

Classifiers assign incoming traffic to forwarding classes (which are mapped to output queues), and also assign a PLP to the incoming traffic. The PLP can be low, medium-high, or high. You can classify traffic

with different PLPs into the same forwarding class to differentiate treatment of traffic within the forwarding class.

In a drop profile map, you can configure a different drop profile for each PLP and associate (map) the drop profiles to a queue scheduler. A scheduler map maps the queue scheduler to a forwarding class (output queue). Traffic classified into the forwarding class uses the drop characteristics defined in the drop profiles that the drop profile map associates with the queue scheduler. The drop profile the traffic uses depends on the PLP that the classifier assigns to the traffic. (You can map different drop profiles to the forwarding class for different PLPs.)

In summary:

- Classifiers assign one of three PLPs (low, medium-high, high) to incoming traffic when classifiers assign traffic to a forwarding class.
- Drop profiles set thresholds for packet drop at different queue fill levels.
- Drop profile maps associate a drop profile with each PLP, and map the drop profiles to schedulers.
- Scheduler maps map schedulers to forwarding classes, and forwarding classes are mapped to output queues. The scheduler mapped to a forwarding class determines the CoS characteristics of the output queue mapped to the forwarding class, including the drop profile mapping.

Buffer Size

Most of the total system buffer space is divided into two buffer pools, shared buffers and dedicated buffers. Shared buffers are a global pool that the ports share dynamically as needed. Dedicated buffers are a reserved portion of the buffer pool that is distributed evenly to all of the ports. Each port receives an equal allocation of dedicated buffer space. The dedicated buffer allocation to ports is not configurable because it is reserved for the ports.

The queue buffers are allocated from the dedicated buffer pool assigned to the port. By default, ports divide their allocation of dedicated buffers among the egress queues in the same proportion as the default scheduler sets the minimum guaranteed transmission rates (`transmit-rate`) for traffic. Only the queues included in the default scheduler receive dedicated buffers.

If you do not use the default configuration, you can explicitly configure the queue buffer size in either of two ways:

- As a percentage—The queue receives the specified percentage of dedicated port buffers when the queue is mapped to the scheduler and the scheduler is mapped to a port.
- As a remainder—After the port services the queues that have an explicit percentage buffer size configuration, the remaining port dedicated buffer space is divided equally among the other queues to which a scheduler is attached. (No default or explicit scheduler means no dedicated buffer

allocation for the queue.) If you configure a scheduler and you do not specify a buffer size as a percentage, *remainder* is the default setting.

NOTE: The total of all of the explicitly configured buffer size percentages for all of the queues on a port cannot exceed 100 percent.

For a complete discussion about queue buffer configuration in the context of ingress and egress port buffer configuration, see [Understanding CoS Buffer Configuration](#).

Explicit Congestion Notification

Explicit congestion notification (ECN) notifies networks about congestion with the goal of reducing packet loss and delay by making the sending device decrease the transmission rate until the congestion clears, without dropping packets. ECN enables end-to-end congestion notification between two endpoints on TCP/IP based networks. ECN is disabled by default.

For more information, see [Understanding CoS Explicit Congestion Notification](#).

Scheduler Maps

A scheduler map associates a forwarding class with a scheduler configuration. After configuring a scheduler, you must include it in a scheduler map, associate the scheduler map with a traffic control profile, and then associate the traffic control profile with an interface and a forwarding class set to implement the configured queue scheduling.

You can associate up to four user-defined scheduler maps with traffic control profiles. For more information, see [Default Schedulers Overview](#).

RELATED DOCUMENTATION

[Understanding Junos CoS Components](#)

[Understanding CoS Priority Group Scheduling](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\)](#)

[Understanding CoS Buffer Configuration](#)

[Understanding CoS Explicit Congestion Notification](#)

[Understanding CoS Scheduling Behavior and Configuration Considerations](#)

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

Configuring CoS Drop Profile Maps

Defining CoS Queue Scheduling Priority

Understanding CoS Priority Group Shaping and Queue Shaping (Maximum Bandwidth)

Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth

Defining CoS Queue Schedulers

Schedulers define the CoS properties of output queues (output queues are mapped to forwarding classes, and classifiers map traffic into forwarding classes based on IEEE 802.1p, DSCP, or MPLS EXP code points). Queue scheduling works with priority group scheduling to create a two-tier hierarchical scheduler. CoS scheduling properties include the amount of interface bandwidth assigned to the queue, the priority of the queue, whether explicit congestion notification (ECN) is enabled on the queue, and the WRED packet drop profiles associated with the queue.

The parameters you configure in a scheduler define the following characteristics for the queues mapped to the scheduler:

- **transmit-rate**—Minimum bandwidth, also known as the *committed information rate (CIR)*, set as a percentage rate or as an absolute value in bits per second. The transmit rate also determines the amount of excess (extra) priority group bandwidth that the queue can share. Extra priority group bandwidth is allocated among the queues in the priority group in proportion to the transmit rate of each queue.

NOTE: Include the preamble bytes and interframe gap (IFG) bytes as well as the data bytes in your bandwidth calculations.

NOTE: You cannot configure a transmit rate for strict-high priority queues. Queues (forwarding classes) with a configured transmit rate cannot be included in a forwarding class set that has strict-high priority queues.

- **shaping-rate**—Maximum bandwidth, also known as the *peak information rate (PIR)*, set as a percentage rate or as an absolute value in bits per second.

NOTE: Include the preamble bytes and interframe gap (IFG) bytes as well as the data bytes in your bandwidth calculations.

- **priority**—One of two bandwidth priorities that queues associated with a scheduler can receive:
 - **low**—The scheduler has low priority.
 - **strict-high**—The scheduler has strict-high priority. You can configure only one queue as a strict-high priority queue. Strict-high priority allocates the scheduled bandwidth to the queue before any other queue receives bandwidth. Other queues receive the bandwidth that remains after the strict-high queue has been serviced.

We recommend that you always apply a shaping rate to strict-high priority queues to prevent them from starving other queues. If you do not apply a shaping rate to limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

- **drop-profile-map**—Drop profile mapping to a loss priority and protocol, to apply WRED to the scheduler and control packet drop for different packet loss priorities during periods of congestion.
- **buffer-size**—Size of the queue buffer as a percentage of the dedicated buffer space on the port, or as a proportional share of the dedicated buffer space on the port that remains after the explicitly configured queues are served.
- **explicit-congestion-notification**—Enables ECN on a best-effort queue. ECN enables end-to-end congestion notification between two ECN-enabled endpoints on TCP/IP based networks. ECN must be enabled on both endpoints and on all of the intermediate devices between the endpoints for ECN to work properly. ECN is disabled by default.

NOTE: Ingress port congestion can occur during periods of egress port congestion if an ingress port forwards traffic to more than one egress port, and at least one of those egress ports experiences congestion. If this occurs, the congested egress port can cause the ingress port to exceed its fair allocation of ingress buffer resources. When the ingress port exceeds its buffer resource allocation, frames are dropped at the ingress. Ingress port frame drop affects not only the congested egress ports, but also all of the egress ports to which the congested ingress port forwards traffic.

If a congested ingress port drops traffic that is destined for one or more uncongested egress ports, configure a weighted random early detection (WRED) drop profile and apply it to the egress queue that is causing the congestion. The drop profile prevents the congested egress queue from affecting egress queues on other ports by dropping frames at the egress instead of causing congestion at the ingress port.

NOTE: Do not configure drop profiles for the fcoe and no-loss forwarding classes. FCoE and other lossless traffic queues require lossless behavior. Use priority-based flow control (PFC) to prevent frame drop on lossless priorities.

OCX Series switches do not support lossless transport or PFC. On OCX Series switches, do not map traffic to the default lossless fcoe and no-loss forwarding classes.

To apply scheduling properties to traffic, map schedulers to forwarding classes using a scheduler map, and then associate the scheduler map with interfaces. (You associate a scheduler map with an interface using a traffic control profile; see [Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#) for an example of the complete hierarchical scheduling process.) Using different scheduler maps, you can map different schedulers to the same traffic (the same forwarding class) on different interfaces, to apply different scheduling to that traffic on different interfaces.

To configure a scheduler using the CLI:

1. Name the scheduler and set the minimum guaranteed bandwidth for the queue:

```
[edit class-of-service]
user@switch# set schedulers scheduler-name transmit-rate (rate | percent
percentage)
```

2. Set the maximum bandwidth for the queue:

```
[edit class-of-service schedulers scheduler-name]
user@switch# set shaping-rate (rate | percent percentage)
```

3. Set the queue priority:

```
[edit class-of-service schedulers scheduler-name]
user@switch# set priority level
```

4. Specify drop profiles for packet loss priorities using a drop profile map:

```
[edit class-of-service schedulers scheduler-name]
user@switch# set drop-profile-map loss-priority (low | medium-high | high) protocol protocol
drop-profile drop-profile-name
```

5. Configure the size of the port dedicated buffer space for the queue:

```
[edit class-of-service schedulers scheduler-name]
user@switch# set buffer-size (percent percent | remainder)
```

6. Enable ECN, if desired (on best-effort traffic only):

```
[edit class-of-service schedulers scheduler-name]
user@switch# set explicit-congestion-notification
```

7. Configure a scheduler map to map the scheduler to a forwarding class, which applies the scheduler's properties to the traffic in that forwarding class:

```
[edit class-of-service]
user@switch# set scheduler-maps scheduler-map-name forwarding-class forwarding-class-name
scheduler scheduler-name
```

8. Assign the scheduler map and its associated schedulers to one or more interfaces using hierarchical scheduling. See [Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#) for a detailed example of hierarchical scheduling.

```
[edit class-of-service]
user@switch# set traffic-control-profiles tcp-name scheduler-map scheduler-map-name
user@switch# set interfaces interface-name forwarding-class-set fc-set-name output-traffic-
control-profile tcp-name
```

RELATED DOCUMENTATION

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#)

[Example: Configuring Queue Schedulers](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth](#)

[Example: Configuring Maximum Output Bandwidth](#)

[Example: Configuring ECN](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\)](#)

[Defining CoS Queue Scheduling Priority](#)

Configuring CoS WRED Drop Profiles

Monitoring CoS Scheduler Maps

Understanding CoS Output Queue Schedulers

Understanding CoS Priority Group Scheduling

Understanding CoS Buffer Configuration

Understanding CoS Explicit Congestion Notification

Example: Configuring Queue Schedulers

IN THIS SECTION

- [Requirements | 200](#)
- [Overview | 200](#)
- [Verification | 203](#)

Schedulers define the CoS properties of output queues (output queues are mapped to forwarding classes, and classifiers map traffic into forwarding classes based on IEEE 802.1p or DSCP code points). Queue scheduling works with priority group scheduling to create a two-tier hierarchical scheduler. CoS scheduling properties include the amount of interface bandwidth assigned to the queue, the priority of the queue, whether explicit congestion notification (ECN) is enabled on the queue, and the WRED packet drop profiles associated with the queue.

Configuring a CoS Scheduler

CLI Quick Configuration

To quickly configure a queue scheduler, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level:

```
[edit class-of-service]
set schedulers be-sched transmit-rate percent 20
set schedulers be-sched shaping-rate percent 40
set schedulers be-sched buffer-size percent 20
```

```

set schedulers be-sched priority low
set schedulers be-sched drop-profile-map loss-priority low protocol any drop-profile be-dp
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set traffic-control-profiles be-tcp scheduler-map be-map
set interfaces xe-0/0/7 forwarding-class-set lan-pg output-traffic-control-profile be-tcp

```

Step-by-Step Procedure

To configure a CoS scheduler:

1. Create scheduler (be-sched) with a minimum guaranteed bandwidth of 2 Gbps, a maximum bandwidth of 4 Gbps, and low priority, and map it to the drop profile be-dp:

```

[edit class-of-service schedulers]
user@switch# set be-sched transmit-rate percent 20
user@switch# set be-sched shaping-rate percent 40
user@switch# set be-sched buffer-size percent 20
user@switch# set be-sched priority low
user@switch# set be-sched drop-profile-map loss-priority low protocol any drop-profile be-
dp

```

NOTE: Because ECN is disabled by default, no ECN configuration is shown.

2. Configure scheduler map (be-map) to associate the scheduler (be-sched) with the forwarding class (best-effort):

```

[edit class-of-service scheduler-maps]
user@switch# set be-map forwarding-class best-effort scheduler be-sched

```

3. Associate the scheduler map be-map with a traffic control profile (be-tcp):

```

[edit class-of-service traffic-control-profiles]
user@switch# set be-tcp scheduler-map be-map

```

4. Associate the traffic control profile `be-tcp` with a forwarding class set (`lan-pg`) and a 10-Gigabit Ethernet interface (`xe-0/0/7`):

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/7 forwarding-class-set lan-pg output-traffic-control-
profile be-tcp
```

Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

Overview

Scheduler parameters define the following characteristics for the queues mapped to the scheduler:

- `transmit-rate`—Minimum bandwidth, also known as the *committed information rate (CIR)*. Each queue mapped to the scheduler receives a minimum of either the configured amount of absolute bandwidth or the configured percentage of bandwidth. The transmit rate also determines the amount of excess (extra) priority group bandwidth that the queue can share. Extra priority group bandwidth is allocated among the queues in the priority group in proportion to the transmit rate of each queue. You cannot configure a transmit rate for strict-high priority queues. Queues (forwarding classes) with a configured transmit rate cannot be included in a forwarding class set that has strict-high priority queues.

NOTE: The `transmit-rate` setting works only if you also configure the `guaranteed-rate` in the traffic control profile that is attached to the forwarding class set to which the queue belongs. If you do not configure the `guaranteed-rate`, the `transmit-rate` does not work. The sum of all queue transmit rates in a forwarding class set should not exceed the traffic control profile guaranteed rate. If you configure transmit rates whose sum exceeds the forwarding class set guaranteed rate, the commit check fails, and the system rejects the configuration.

NOTE: Include the preamble bytes and interframe gap bytes as well as the data bytes in your bandwidth calculations.

NOTE: You cannot configure a transmit rate for strict-high priority queues. Queues (forwarding classes) with a configured transmit rate cannot be included in a forwarding class set that has strict-high priority queues.

- **shaping-rate**—Maximum bandwidth, also known as the *peak information rate (PIR)*. Each queue receives a maximum of the configured amount of absolute bandwidth or the configured percentage of bandwidth, even if more bandwidth is available.

NOTE: Include the preamble bytes and interframe gap bytes as well as the data bytes in your bandwidth calculations.

- **priority**—One of two bandwidth priorities that queues associated with a scheduler can receive:
 - **low**—The scheduler has low priority.
 - **strict-high**—The scheduler has strict-high priority. You can configure only one queue as a strict-high priority queue. Strict-high priority allocates the scheduled bandwidth to the queue before any other queue receives bandwidth. Other queues receive the bandwidth that remains after the strict-high queue has been serviced.

We recommend that you always apply a shaping rate to strict-high priority queues to prevent them from starving other queues. If you do not apply a shaping rate to limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

- **drop-profile-map**—Mapping of a drop profile to a loss priority and protocol to apply WRED to the scheduler.
- **buffer-size**—Size of the queue buffer as a percentage of the dedicated buffer space on the port, or as a proportional share of the dedicated buffer space on the port that remains after the explicitly configured queues are served.
- **explicit-congestion-notification**—Enables ECN on a best-effort queue. ECN enables end-to-end congestion notification between two ECN-enabled endpoints on TCP/IP based networks. ECN must be enabled on both endpoints and on all of the intermediate devices between the endpoints for ECN to work properly. ECN is disabled by default.

NOTE: Ingress port congestion can occur during periods of egress port congestion if an ingress port forwards traffic to more than one egress port, and at least one of those egress ports experiences congestion. If this occurs, the congested egress port can cause the ingress port to

exceed its fair allocation of ingress buffer resources. When the ingress port exceeds its buffer resource allocation, frames are dropped at the ingress. Ingress port frame drop affects not only the congested egress ports, but also all of the egress ports to which the congested ingress port forwards traffic.

If a congested ingress port drops traffic that is destined for one or more uncongested egress ports, configure a weighted random early detection (WRED) drop profile and apply it to the egress queue that is causing the congestion. The drop profile prevents the congested egress queue from affecting egress queues on other ports by dropping frames at the egress instead of causing congestion at the ingress port.

NOTE: Do not configure drop profiles for the fcoe and no-loss forwarding classes. FCoE and other lossless traffic queues require lossless behavior. Use priority-based flow control (PFC) to prevent frame drop on lossless priorities.

OCX Series switches do not support lossless transport or PFC. On OCX Series switches, do not map traffic to the default lossless fcoe and no-loss forwarding classes.

Scheduler maps associate schedulers with forwarding classes (queues). After defining schedulers and mapping them to queues in a scheduler map, to configure hardware queue scheduling (hierarchical port scheduling) you:

1. Associate a scheduler map with a traffic control profile (a traffic control profile schedules resources for a group of forwarding classes, called a *forwarding class set* or *priority group*).
2. Attach a forwarding class and a traffic control profile to an interface.

Example: Configuring CoS Hierarchical Port Scheduling (ETS) provides a complete example of hierarchical scheduling.

You can associate up to four user-defined scheduler maps with forwarding class sets.

This process configures the bandwidth properties and WRED characteristics that you map to forwarding classes (and thus to output queues) in a scheduler map. The traffic control profile uses the scheduler CoS properties to determine the resources that should be allocated to the individual output queues from the total resources available to the priority group.

[Table 57 on page 203](#) shows the configuration components for this example.

Table 57: Components of the Queue Scheduler Configuration Example

Component	Settings
Hardware	QFX3500 switch
Scheduler	Name: be-sched Transmit rate: 20% Shaping rate: 40% Buffer size: 20% Priority: low Drop profile: be-dp ECN: disable (default)
Scheduler map	Name: be-map Forwarding class to associate with the be-sched scheduler: best-effort
Traffic control profile	Name: be-tcp NOTE: This topic does not describe how to define a traffic control profile.
Forwarding class set	Name: lan-pg

Verification

IN THIS SECTION

- [Verifying the Scheduler Configuration | 204](#)
- [Verifying the Scheduler Map Configuration | 204](#)
- [Verifying That the Scheduler Is Associated with the Interface | 205](#)

To verify that the queue scheduler has been created and is mapped to the correct interfaces, perform these tasks:

Verifying the Scheduler Configuration

Purpose

Verify that the queue scheduler `be-sched` has been created with a minimum guaranteed bandwidth of 2 Gbps, a maximum bandwidth of 4 Gbps, the priority set to `low`, and the drop profile `be-dp`.

Action

Display the scheduler using the operational mode command `show configuration class-of-service schedulers be-sched`:

```
user@switch> show configuration class-of-service schedulers be-sched
transmit-rate percent 20;
shaping-rate percent 40;
buffer-size percent 20;
priority low;
drop-profile-map loss-priority low protocol any drop-profile be-dp;
```

Verifying the Scheduler Map Configuration

Purpose

Verify that the scheduler map `be-map` has been created and associates the forwarding class `best-effort` with the scheduler `be-sched`, and also that the scheduler map is attached to the traffic control profile `be-tcp`.

Action

Display the scheduler map using the operational mode command `show configuration class-of-service scheduler-maps be-map`:

```
user@switch> show configuration class-of-service scheduler-maps be-map
forwarding-class best-effort scheduler be-sched;
```

Display the traffic control profile to verify that the scheduler map `be-map` is attached using the operational mode command `show configuration class-of-service traffic-control-profiles be-tcp scheduler-map`:

```
user@switch> show configuration class-of-service traffic-control-profiles be-tcp scheduler-map
scheduler-map be-map;
```

NOTE: This topic does not describe how to configure a traffic control profile or its allocation of port bandwidth. Using a traffic control profile to configure the port resource allocation to the priority group is necessary to implement hierarchical scheduling.

Verifying That the Scheduler Is Associated with the Interface

Purpose

Verify that the forwarding class set (`lan-pg`) and the traffic control profile (`be-tcp`) that are associated with the queue scheduler are attached to the interface `xe-0/0/7`.

Action

List the interface using the operational mode command `show configuration class-of-service interfaces xe-0/0/7`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/7
forwarding-class-set {
    lan-pg {
        output-traffic-control-profile be-tcp;
    }
}
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Minimum Guaranteed Output Bandwidth

Example: Configuring Maximum Output Bandwidth

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Example: Configuring WRED Drop Profiles

Example: Configuring ECN

Defining CoS Queue Schedulers

Monitoring CoS Scheduler Maps

Understanding CoS Output Queue Schedulers

Understanding CoS Hierarchical Port Scheduling (ETS)

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

Understanding CoS Buffer Configuration

Defining CoS Queue Scheduling Priority

You can configure the scheduling priority of individual queues by specifying the priority in a scheduler, and then associating the scheduler with a queue by using a scheduler map. On QFX5100, QFX5200, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, queues can have one of two bandwidth scheduling priorities, strict-high priority or low priority. On QFX10000 Series switches, queues can also be configured as high priority.

NOTE: By default, all queues are low priority queues.

The switch services low priority queues after servicing any queue that has strict-high priority traffic or high priority traffic. Strict-high priority queues receive preferential treatment over all other queues and receive all of their configured bandwidth before other queues are serviced. Low-priority queues do not transmit traffic until strict-high priority queues are empty, and receive the bandwidth that remains after the strict-high queues have been serviced. High priority queues receive preference over low priority queues.

Different switches handle traffic configured as strict-high priority traffic in different ways:

- QFX5100, QFX5200, QFX3500, QFX3600, and EX4600 switches, and QFabric systems—You can configure only one queue as a strict-high priority queue.

On these switches, we recommend that you always apply a shaping rate to strict-high priority queues to prevent them from starving other queues. If you do not apply a shaping rate to limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

- QFX10000 switches—You can configure as many queues as you want as strict-high priority. However, keep in mind that too much strict-high priority traffic can starve low priority queues on the port.

NOTE: We strongly recommend that you configure a transmit rate on all strict-high priority queues to limit the amount of traffic the switch treats as strict-high priority traffic and prevent strict-high priority queues from starving other queues on the port. This is especially important if you configure more than one strict-high priority queue on a port. If you do not configure a transmit rate to limit the amount of bandwidth strict-high priority queues can use, then the strict-high priority queues can use all of the available port bandwidth and starve other queues on the port.

The switch treats traffic in excess of the transmit rate as best-effort traffic that receives bandwidth from the leftover (excess) port bandwidth pool. On strict-high priority queues, all traffic that exceeds the transmit rate shares in the port excess bandwidth pool based on the strict-high priority excess bandwidth sharing weight of “1”, which is not configurable. The actual amount of extra bandwidth that traffic exceeding the transmit rate receives depends on how many other queues consume excess bandwidth and the excess rates of those queues.

- To configure queue priority using the CLI:

```
[edit class-of-service]
```

```
user@switch# set schedulers scheduler-name priority level
```

RELATED DOCUMENTATION

Example: Configuring Queue Scheduling Priority

Monitoring CoS Scheduler Maps

Example: Configuring Queue Scheduling Priority

IN THIS SECTION

- Requirements | 209
- Overview | 209
- Verification | 211

You can configure the bandwidth scheduling priority of individual queues by specifying the priority in a scheduler, and then using a scheduler map to associate the scheduler with a queue.

Configuring Queue Scheduling Priority

CLI Quick Configuration

To quickly configure queue scheduling priority, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level:

```
[edit class-of-service]
set schedulers fcoe-sched priority low
set schedulers nl-sched priority low
set scheduler-maps schedmap1 forwarding-class fcoe scheduler fcoe-sched
set scheduler-maps schedmap1 forwarding-class no-loss scheduler nl-sched
```

Step-by-Step Procedure

To configure queue priority using the CLI:

1. Create the FCoE scheduler with low priority:

```
[edit class-of-service]
user@switch# set schedulers fcoe-sched priority low
```

2. Create the no-loss scheduler with low priority:

```
[edit class-of-service]
user@switch# set schedulers nl-sched priority low
```

3. Associate the schedulers with the desired queues in the scheduler map:

```
[edit class-of-service]
user@switch# set scheduler-maps schedmap1 forwarding-class fcoe scheduler fcoe-sched
user@switch# set scheduler-maps schedmap1 forwarding-class no-loss scheduler nl-sched
```

Requirements

This example uses the following hardware and software components:

- One switch.
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series.

Overview

Queues can have one of several bandwidth priorities:

- **strict-high**—Strict-high priority allocates bandwidth to the queue before any other queue receives bandwidth. Other queues receive the bandwidth that remains after the strict-high queue has been serviced. On QFX10000 switches, you can configure as many queues as you want as strict-high priority queues. On QFX5200, QFX3500, and QFX3600 switches and on QFabric systems, you can configure only one queue as a strict-high queue. On QFX5100 and EX4600 switches, you can configure only one forwarding-class-set (priority group) as strict-high priority. All queues which are part of that strict-high forwarding class set then act as strict-high queues.

NOTE: On QFX5200 switches, it is not possible to support multiple queues with strict-high priority because QFX5200 doesn't support flexible hierarchical scheduling. When multiple strict-high priority queues are configured, all of those queues are treated as strict-high priority but the higher number queue among them is given highest priority.

On QFX10000 switches, if you configure strict-high priority queues on a port, we strongly recommend that you configure a transmit rate on those queues. The transmit rate sets the amount of traffic that the switch forwards as strict-high priority; traffic in excess of the transmit rate is treated as best-effort traffic that receives the queue excess rate. Even if you configure only one strict-high priority queue, we strongly recommend that you configure a transmit rate the queue to prevent it from starving other queues. If you do not configure a transmit rate to limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

On QFX5200, QFX5100, QFX3500, QFX3600, and EX4600 switches and on QFabric systems, we recommend that you always apply a shaping rate to strict-high priority queues to prevent them from starving other queues. If you do not apply a shaping rate to limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

NOTE: On switches that support enhanced transmission selection (ETS) hierarchical scheduling, if you use ETS and you configure a strict-high priority queue, you must create a forwarding class set that is dedicated only to strict-high priority traffic. Only one forwarding class set can contain a strict-high priority queue. Queues that are not strict-high priority cannot belong to the same forwarding class set as strict-high priority queues.

On switches that use different output queues for unicast and multideestination traffic, the multideestination forwarding class set cannot contain strict-high priority queues.

- `high` (QFX10000 Series switches only)—High priority. Traffic with high priority is serviced after any queue that has a strict-high priority, and before queues with low priority.
- `low`—Low priority. Traffic with low priority is serviced after any queue that has a strict-high priority.

NOTE: By default, all queues are low priority queues.

Table 58 on page 210 shows the configuration components for this example.

This example describes how to set the queue priority for two forwarding classes (queues) named `fcoe` and `no-loss`. Both queues have a priority of `low`. The scheduler for the `fcoe` queue is named `fcoe-sched` and the scheduler for the `no-loss` queue is named `n1-sched`. One scheduler map, `schedmap1`, associates the schedulers to the queues.

Table 58: Components of the Queue Scheduler Priority Configuration Example

Component	Settings
Hardware	One switch
Schedulers	<code>fcoe-sched</code> for FCoE traffic <code>n1-sched</code> for no-loss traffic
Priority	<code>low</code> for FCoE traffic <code>low</code> for no-loss traffic

Table 58: Components of the Queue Scheduler Priority Configuration Example *(Continued)*

Component	Settings
Scheduler map	<p>schedmap1:</p> <p>FCoE mapping: scheduler fcoe-sched to forwarding class fcoe</p> <p>No-loss mapping: scheduler nl-sched to forwarding class no-loss</p>

NOTE: OCX Series switches do not support lossless transport. On OCX Series switches, the default DSCP classifier does not map traffic to the default fcoe and no-loss forwarding classes. On an OCX Series switch, you could use this example by substituting other forwarding classes (for example, best-effort or network-control) for the fcoe and no-loss forwarding classes, and naming the schedulers appropriately. The active forwarding classes (best-effort, network-control, and mcast) share the unused bandwidth assigned to the fcoe and no-loss forwarding classes.

Verification

IN THIS SECTION

- [Verifying the Queue Scheduling Priority | 211](#)
- [Verifying the Scheduler-to-Forwarding-Class Mapping | 212](#)

To verify that you configured the queue scheduling priority for bandwidth and mapped the schedulers to the correct forwarding classes, perform these tasks:

Verifying the Queue Scheduling Priority

Purpose

Verify that you configured the queue schedulers fcoe-sched and nl-sched with low queue scheduling priority.

Action

Display the fcoe-sched scheduler priority configuration using the operational mode command `show configuration class-of-service schedulers fcoe-sched priority`:

```
user@switch> show configuration class-of-service schedulers fcoe-sched priority
priority low;
```

Display the nl-sched scheduler priority configuration using the operational mode command `show configuration class-of-service schedulers nl-sched priority`:

```
user@switch> show configuration class-of-service schedulers nl-sched priority
priority low;
```

Verifying the Scheduler-to-Forwarding-Class Mapping

Purpose

Verify that you configured the scheduler map `schedmap1` to map scheduler `fcoe-sched` to forwarding class `fcoe` and schedule `nl-sched` to forwarding class `no-loss`.

Action

Display the scheduler map `schedmap1` using the operational mode command `show configuration class-of-service scheduler-maps schedmap1`:

```
user@switch> show configuration class-of-service scheduler-maps schedmap1
forwarding-class fcoe scheduler fcoe-sched;
forwarding-class no-loss scheduler nl-sched;
```

RELATED DOCUMENTATION

Defining CoS Queue Scheduling Priority

Monitoring CoS Scheduler Maps

Understanding CoS Traffic Control Profiles

A traffic control profile defines the output bandwidth and scheduling characteristics of forwarding class sets (priority groups). The forwarding classes (which are mapped to output queues) that belong to a forwarding class set (fc-set) share the bandwidth that you assign to the fc-set in the traffic control profile.

This two-tier hierarchical scheduling architecture provides flexibility in allocating resources among forwarding classes, and also:

- Assigns a portion of port bandwidth to an fc-set. You define the port resources for the fc-set in a traffic control profile.
- Allocates fc-set bandwidth among the forwarding classes (queues) that belong to the fc-set. A scheduler map attached to the traffic control profile defines the amount of the fc-set's resources that each forwarding class can use.

Attaching an fc-set and a traffic control profile to a port defines the hierarchical scheduling properties of the group and the forwarding classes that belong to the group.

The ability to create fc-sets supports enhanced transmission selection (ETS), which is described in IEEE 802.1Qaz. When an fc-set does not use its allocated port bandwidth, ETS shares the excess port bandwidth among other fc-sets on the port in proportion to their guaranteed minimum bandwidth (guaranteed rate). This utilizes the port bandwidth better than scheduling schemes that reserve bandwidth for groups even if that bandwidth is not used. ETS shares unused port bandwidth, so traffic groups that need extra bandwidth can use it if the bandwidth is available, while preserving the ability to specify the minimum guaranteed bandwidth for traffic groups.

Traffic control profiles define the following CoS properties for fc-sets:

- Minimum guaranteed bandwidth—Also known as the *committed information rate (CIR)*. This is the minimum amount of port bandwidth the priority group receives. Priorities in the priority group receive their minimum guaranteed bandwidth as a portion of the priority group's minimum guaranteed bandwidth. The guaranteed-rate statement defines the minimum guaranteed bandwidth.

NOTE: You cannot apply a traffic control profile with a minimum guaranteed bandwidth to a priority group that includes strict-high priority queues.

- Shared excess (extra) bandwidth—When the priority groups on a port do not consume the full amount of bandwidth allocated to them or there is unallocated link bandwidth available, priority groups can contend for that extra bandwidth if they need it. Priorities in the priority group contend for extra bandwidth as a portion of the priority group's extra bandwidth. The amount of extra

bandwidth for which a priority group can contend is proportional to the priority group's guaranteed minimum bandwidth (guaranteed rate).

- **Maximum bandwidth**—Also known as *peak information rate (PIR)*. This is the maximum amount of port bandwidth the priority group receives. Priorities in the priority group receive their maximum bandwidth as a portion of the priority group's maximum bandwidth. The `shaping-rate` statement defines the maximum bandwidth.
- **Queue scheduling**—Each traffic control profile includes a scheduler map. The scheduler map maps forwarding classes (priorities) to schedulers to define the scheduling characteristics of the individual forwarding classes in the fc-set. The resources scheduled for each forwarding class represent portions of the resources that the traffic control profile schedules for the entire fc-set, not portions of the total link bandwidth. The `scheduler-maps` statement defines the mapping of forwarding classes to schedulers.

RELATED DOCUMENTATION

Understanding CoS Hierarchical Port Scheduling (ETS)

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Defining CoS Traffic Control Profiles (Priority Group Scheduling)

Understanding CoS Priority Group Scheduling

IN THIS SECTION

- [Priority Group Scheduling Components | 215](#)
- [Default Traffic Control Profile | 216](#)
- [Guaranteed Rate \(Minimum Guaranteed Bandwidth\) | 216](#)
- [Sharing Extra Bandwidth | 217](#)
- [Shaping Rate \(Maximum Bandwidth\) | 217](#)
- [Scheduler Maps | 217](#)

Priority group scheduling defines the class-of-service (CoS) properties of a group of output queues (priorities). Priority group scheduling works with output queue scheduling to create a two-tier

hierarchical scheduler. The hierarchical scheduler allocates bandwidth to a group of queues (a priority group, called a forwarding class set in Junos OS configuration). Queue scheduling determines the portion of the priority group bandwidth that the particular queue can use.

You configure priority group scheduling in a traffic control profile and then associate the traffic control profile with a forwarding class set and an interface. You attach a scheduler map to the traffic control profile to specify the queue scheduling characteristics.

NOTE: When you configure bandwidth for a queue or a priority group, the switch considers only the data as the configured bandwidth. The switch does not account for the bandwidth consumed by the preamble and the interframe gap (IFG). Therefore, when you calculate and configure the bandwidth requirements for a queue or for a priority group, consider the preamble and the IFG as well as the data in the calculations.

Priority Group Scheduling Components

Table 59 on page 215 provides a quick reference to the traffic control profile components you can configure to determine the bandwidth properties of priority groups, and Table 60 on page 216 provides a quick reference to some related scheduling configuration components.

Table 59: Priority Group Scheduler Components

Traffic Control Profile Component	Description
Guaranteed rate	Sets the minimum guaranteed port bandwidth for the priority group. Extra port bandwidth is shared among priority groups in proportion to the guaranteed rate of each priority group on the port.
Shaping rate	Sets the maximum port bandwidth the priority group can consume.
Scheduler map	Maps schedulers to queues (forwarding classes, also called priorities). This determines the portion of the priority group bandwidth that a queue receives.

Table 60: Other Scheduling Components

Other Scheduling Components	Description
Forwarding class	Maps traffic to a queue (priority).
Forwarding class set	Name of a priority group. You map forwarding classes to priority groups. A forwarding class set consists of one or more forwarding classes.
Scheduler	Sets the bandwidth and scheduling priority of individual queues (forwarding classes).

Default Traffic Control Profile

There is no default traffic control profile.

Guaranteed Rate (Minimum Guaranteed Bandwidth)

The guaranteed rate determines the minimum guaranteed bandwidth for each priority group. It also determines how much excess (extra) port bandwidth the priority group can share; each priority group shares extra port bandwidth in proportion to its guaranteed rate. You specify the rate in bits per second as a fixed value such as 3 Mbps or as a percentage of the total port bandwidth.

The minimum transmission bandwidth can exceed the configured rate if additional bandwidth is available from other priority groups on the port. In case of congestion, the configured guaranteed rate is guaranteed for the priority group. This property enables you to ensure that each priority group receives the amount of bandwidth appropriate to its level of service.

NOTE: Configuring the minimum guaranteed bandwidth (transmit rate) for a forwarding class does not work unless you also configure the minimum guaranteed bandwidth (guaranteed rate) for the forwarding class set in the traffic control profile.

Additionally, the sum of the transmit rates of the queues in a forwarding class set should not exceed the guaranteed rate for the forwarding class set. (You cannot guarantee a minimum bandwidth for the queues that is greater than the minimum bandwidth guaranteed for the entire set of queues.)

You cannot configure a guaranteed rate for forwarding class sets that include strict-high priority queues.

Sharing Extra Bandwidth

Extra bandwidth is available to priority groups when the priority groups do not use the full amount of available port bandwidth. This extra port bandwidth is shared among the priority groups based on the minimum guaranteed bandwidth of each priority group.

For example, Port A has three priority groups: fc-set-1, fc-set-2, and fc-set-3. Fc-set-1 has a guaranteed rate of 2 Gbps, fc-set-2 has a guaranteed rate of 2 Gbps, and fc-set-3 has a guaranteed rate of 4 Gbps. After servicing the minimum guaranteed bandwidth of these priority groups, the port has an extra 2 Gbps of available bandwidth, and all three priority groups have still have packets to forward. The priority groups receive the extra bandwidth in proportion to their guaranteed rates, so fc-set-1 receives an extra 500 Mbps, fc-set-2 receives an extra 500 Mbps, and fc-set-3 receives an extra 1 Gbps.

Shaping Rate (Maximum Bandwidth)

The shaping rate determines the maximum bandwidth the priority group can consume. You specify the rate in bits per second as a fixed value such as 5 Mbps or as a percentage of the total port bandwidth.

The maximum bandwidth for a priority group depends on the total bandwidth available on the port and how much bandwidth the other priority groups on the port consume.

Scheduler Maps

A scheduler map maps schedulers to queues. When you associate a scheduler map with a traffic control profile, then associate the traffic control profile with an interface and a forwarding class set, the scheduling defined by the scheduler map determines the portion of the priority group resources that each individual queue can use.

You can associate up to four user-defined scheduler maps with traffic control profiles.

RELATED DOCUMENTATION

Understanding Junos CoS Components

Understanding CoS Output Queue Schedulers

Understanding CoS Hierarchical Port Scheduling (ETS)

Understanding CoS Scheduling Behavior and Configuration Considerations

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Minimum Guaranteed Output Bandwidth

Example: Configuring Maximum Output Bandwidth

Example: Configuring Queue Schedulers

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Example: Configuring WRED Drop Profiles

Example: Configuring Drop Profile Maps

Defining CoS Traffic Control Profiles (Priority Group Scheduling)

A traffic control profile defines the output bandwidth and scheduling characteristics of forwarding class sets (priority groups). The forwarding classes (which are mapped to output queues) contained in a forwarding class set (fc-set) share the bandwidth resources that you configure in the traffic control profile. A scheduler map associates forwarding classes with schedulers to define how the individual forwarding classes that belong to an fc-set share the bandwidth allocated to that fc-set.

The parameters you configure in a traffic control profile define the following characteristics for the fc-set:

- **guaranteed-rate**—Minimum bandwidth, also known as the *committed information rate (CIR)*. The guaranteed rate also determines the amount of excess (extra) port bandwidth that the fc-set can share. Extra port bandwidth is allocated among the fc-sets on a port in proportion to the guaranteed rate of each fc-set.

NOTE: You cannot configure a guaranteed rate for a fc-set that includes strict-high priority queues. If the traffic control profile is for an fc-set that contains strict-high priority queues, do not configure a guaranteed rate.

- **shaping-rate**—Maximum bandwidth, also known as the *peak information rate (PIR)*.
- **scheduler-map**—Bandwidth and scheduling characteristics for the queues, defined by mapping forwarding classes to schedulers. (The queue scheduling characteristics represent amounts or percentages of the fc-set bandwidth, not the amounts or percentages of total link bandwidth.)

NOTE: Because a port can have more than one fc-set, when you assign resources to an fc-set, keep in mind that the total port bandwidth must serve all of the queues associated with that port.

To configure a traffic control profile using the CLI:

1. Name the traffic control profile and define the minimum guaranteed bandwidth for the fc-set:

```
[edit class-of-service ]
user@switch# set traffic-control-profiles traffic-control-profile-name guaranteed-rate (rate
| percent percentage)
```

2. Define the maximum bandwidth for the fc-set:

```
[edit class-of-service traffic-control-profiles traffic-control-profile-name]
user@switch# set shaping-rate (rate | percent percentage)
```

3. Attach a scheduler map to the traffic control profile:

```
[edit class-of-service traffic-control-profiles ]
user@switch# set scheduler-map scheduler-map-name
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Example: Configuring Minimum Guaranteed Output Bandwidth

Example: Configuring Maximum Output Bandwidth

Defining CoS Queue Schedulers

Understanding CoS Traffic Control Profiles

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

IN THIS SECTION

- [Requirements | 221](#)
- [Overview | 221](#)
- [Verification | 222](#)

A traffic control profile defines the output bandwidth and scheduling characteristics of forwarding class sets (priority groups). The forwarding classes (queues) mapped to a forwarding class set share the bandwidth resources that you configure in the traffic control profile. A scheduler map associates forwarding classes with schedulers to define how the individual queues in a forwarding class set share the bandwidth allocated to that forwarding class set.

Configuring a Traffic Control Profile

Step-by-Step Procedure

This example describes how to configure a traffic control profile named `san-tcp` with a scheduler map named `san-map1` and allocate to it a minimum bandwidth of 4 Gbps and a maximum bandwidth of 8 Gbps:

1. Create the traffic control profile and set the `guaranteed-rate` (minimum guaranteed bandwidth) to 4g:

```
[edit class-of-service]
user@switch# set traffic-control-profiles san-tcp guaranteed-rate 4g
```

2. Set the `shaping-rate` (maximum guaranteed bandwidth) to 8g:

```
[edit class-of-service]
user@switch# set traffic-control-profiles san-tcp shaping-rate 8g
```

3. Associate the scheduler map `san-map1` with the traffic control profile:

```
[edit class-of-service]
user@switch# set traffic-control-profiles san-tcp scheduler-map san-map1
```

Requirements

This example uses the following hardware and software components:

- A Juniper Networks QFX3500 Switch
- Junos OS Release 11.1 or later for the QFX Series

Overview

The parameters you configure in a traffic control profile define the following characteristics for the priority group:

- **guaranteed-rate**—Minimum bandwidth, also known as the *committed information rate (CIR)*. Each fc-set receives a minimum of either the configured amount of absolute bandwidth or the configured percentage of bandwidth. The guaranteed rate also determines the amount of excess (extra) port bandwidth that the fc-set can share. Extra port bandwidth is allocated among the fc-sets on a port in proportion to the guaranteed rate of each fc-set.

NOTE: In order for the *transmit-rate* option (minimum bandwidth for a queue that you set using scheduler configuration) to work properly, you must configure the **guaranteed-rate** for the fc-set. If an fc-set does not have a guaranteed minimum bandwidth, the forwarding classes that belong to the fc-set cannot have a guaranteed minimum bandwidth.

NOTE: Include the preamble bytes and interframe gap bytes as well as the data bytes in your bandwidth calculations.

- **shaping-rate**—Maximum bandwidth, also known as the *peak information rate (PIR)*. Each fc-set receives a maximum of the configured amount of absolute bandwidth or the configured percentage of bandwidth, even if more bandwidth is available.

NOTE: Include the preamble bytes and interframe gap bytes as well as the data bytes in your bandwidth calculations.

- **scheduler-map**—Bandwidth and scheduling characteristics for the queues, defined by mapping forwarding classes to schedulers. (The queue scheduling characteristics represent amounts or percentages of the fc-set bandwidth, not the amounts or percentages of total link bandwidth.)

NOTE: Because a port can have more than one fc-set, when you assign resources to an fc-set, keep in mind that the total port bandwidth must serve all of the queues associated with that port.

For example, if you map three fc-sets to a 10-Gigabit Ethernet port, the queues associated with all three of the fc-sets share the 10-Gbps bandwidth as defined by the traffic control profiles. Therefore, the total combined guaranteed-rate value of the three fc-sets should not exceed 10 Gbps. If you configure guaranteed rates whose sum exceeds the port bandwidth, the system sends a syslog message to notify you that the configuration is not valid. However, the system does not perform a commit check. If you commit a configuration in which the sum of the guaranteed rates exceeds the port bandwidth, the hierarchical scheduler behaves unpredictably.

The sum of the forwarding class (queue) transmit rates cannot exceed the total guaranteed-rate of the fc-set to which the forwarding classes belong. If you configure transmit rates whose sum exceeds the fc-set guaranteed rate, the commit check fails and the system rejects the configuration.

If you configure the guaranteed-rate of an fc-set as a percentage, configure all of the transmit rates associated with that fc-set as percentages. In this case, if any of the transmit rates are configured as absolute values instead of percentages, the configuration is not valid and the system sends a syslog message.

Verification

IN THIS SECTION

- [Verifying the Traffic Control Profile Configuration | 222](#)

Verifying the Traffic Control Profile Configuration

Purpose

Verify that you created the traffic control profile `san-tcp` with a minimum guaranteed bandwidth of 4 Gbps, a maximum bandwidth of 8 Gbps, and the scheduler map `san-map1`.

Action

List the traffic control profile using the operational mode command `show configuration class-of-service traffic-control-profiles san-tcp`:

```
user@switch> show configuration class-of-service traffic-control-profiles san-tcp
scheduler-map san-map1;
shaping-rate percent 8g;
guaranteed-rate 4g;
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Minimum Guaranteed Output Bandwidth

Example: Configuring Maximum Output Bandwidth

Example: Configuring Queue Schedulers

Defining CoS Traffic Control Profiles (Priority Group Scheduling)

Understanding CoS Traffic Control Profiles

Understanding CoS Hierarchical Port Scheduling (ETS)

Understanding CoS Hierarchical Port Scheduling (ETS)

IN THIS SECTION

- [Hierarchical Scheduling Tiers | 224](#)
- [Hierarchical Scheduling and ETS | 225](#)
- [ETS Advertisement in DCBX | 227](#)
- [Hierarchical Scheduling Process | 227](#)
- [Strict-High Priority Queues and Hierarchical Scheduling | 229](#)
- [Default Hierarchical Scheduling | 229](#)

Scheduling defines the class-of-service (CoS) properties of output queues. Output queues are mapped to forwarding classes. CoS scheduler properties include the amount of interface bandwidth assigned to the queue, the queue priority, and the drop profiles associated with the queue.

Hierarchical port scheduling is a two-tier process that provides better port bandwidth utilization and greater flexibility to allocate resources to queues (forwarding classes) and to groups of queues (forwarding class sets). Hierarchical scheduling includes the Junos OS implementation of enhanced transmission selection (ETS), as described in IEEE 802.1Qaz.



Video: [What is Enhanced Transmission Selection?](#)

This topic describes:

Hierarchical Scheduling Tiers

The two tiers used in hierarchical scheduling are priorities and priority groups, as shown in [Table 61 on page 224](#).

Table 61: Hierarchical Scheduling Tiers

Junos OS Configuration Construct	Equivalent ETS Construct	Description
Forwarding class	Priority	<p>Think about priorities (forwarding classes) as output queues. You map forwarding classes to queues, so each forwarding class represents an output queue.</p> <p>When you use a classifier to map a forwarding class to an IEEE 802.1p code point, the code point identifies that traffic's priority for priority-based flow control (PFC). Thus the forwarding class, the queue mapped to the forwarding class, and the priority (code point) mapped to the forwarding class all identify the same traffic.</p>
Forwarding class set	Priority group	<p>Priority groups (forwarding class sets) are groups of priorities (forwarding classes). Forwarding class membership in a forwarding class set defines the priority group to which each priority belongs.</p> <p>You can configure up to three unicast priority groups and one multicast priority group.</p>

You apply scheduling properties to each hierarchical scheduling tier as described in the next section.

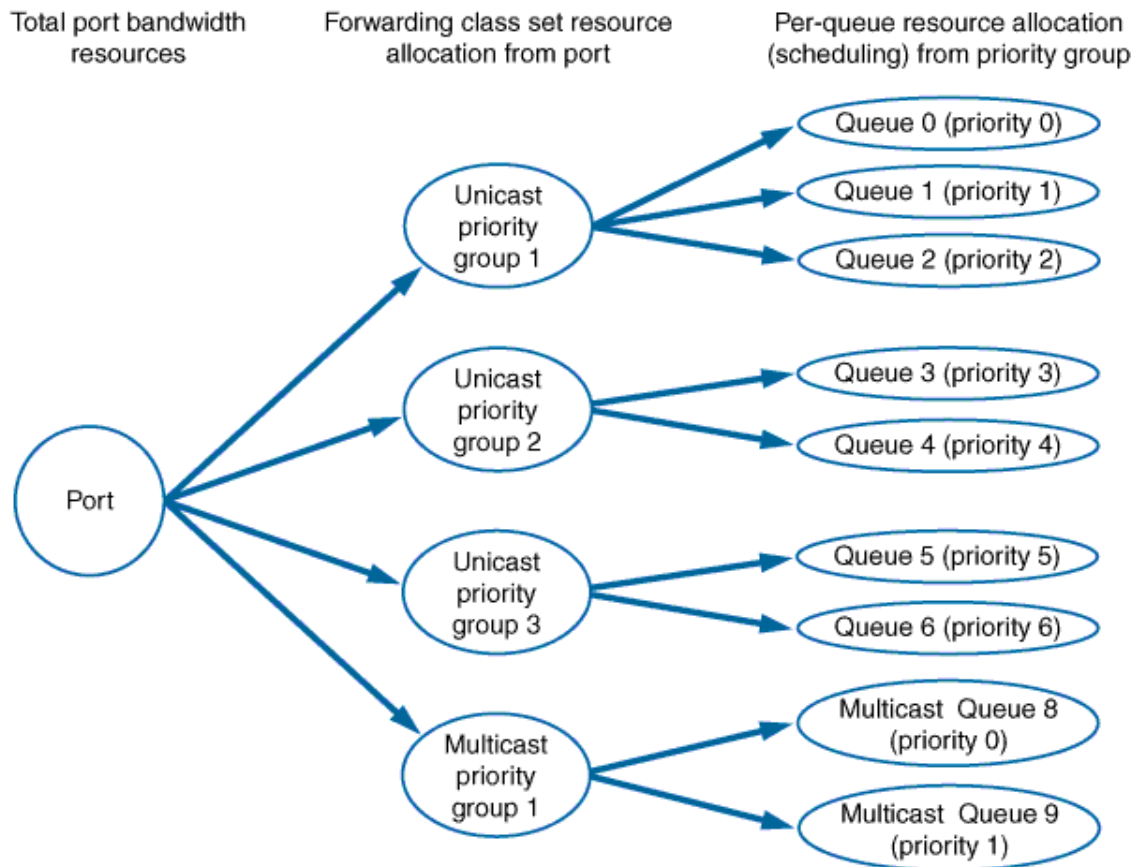
NOTE: If you explicitly configure one or more priority groups on an interface, any priority (forwarding class) that is not assigned to a priority group (forwarding class set) on that interface is assigned to an automatically generated default priority group and receives *no bandwidth*. This means that if you configure hierarchical scheduling on an interface, every forwarding class that you want to forward traffic on that interface must belong to a forwarding class set.

Hierarchical Scheduling and ETS

Two-tier hierarchical scheduling manages bandwidth efficiently by enabling you to define the CoS properties for each priority group and for each priority. The first tier of the hierarchical scheduler allocates port bandwidth to a priority group. The second tier of the hierarchical scheduler determines the portion of the priority group bandwidth that a priority (queue) can use.

The CoS properties of a priority group define the amount of port bandwidth resources available to the queues in that priority group. The CoS properties you configure for each queue specify the amount of the bandwidth available to the queue from the bandwidth allocated to the priority group. [Figure 5 on page 226](#) shows the relationship of port resource allocation to priority groups, and priority group resource allocation to queues (priorities).

Figure 5: Hierarchical Scheduling Tiers



g040722

If a queue (priority) does not use its allocated bandwidth, ETS shares the unused bandwidth among the other queues in the priority group in proportion to the minimum guaranteed rate (transmit rate) scheduled for each queue. If a priority group does not use its allocated bandwidth, ETS shares the unused bandwidth among the priority groups on the port in proportion to the minimum guaranteed rate (guaranteed rate) scheduled for each priority group.

In this way, ETS improves link bandwidth utilization, and it provides each queue and each priority group with the maximum available bandwidth. For example, priorities that consist of bursty traffic can share bandwidth during periods of low traffic transmission, instead of reserving their entire bandwidth allocation when traffic loads are light.

NOTE: The available link bandwidth is the bandwidth remaining after servicing strict-high priority flows. Strict-high priority takes precedence over all other traffic. We recommend that you configure a *shaping-rate* (*transmit-rate* on QFX10000 switches) to limit the maximum amount of bandwidth that a strict-high priority forwarding class can use to prevent starving other queues.

ETS Advertisement in DCBX

When you configure hierarchical scheduling on a port, Data Center Bridging Capability Exchange protocol (DCBX) advertises:

- Each priority group
- The priorities in each priority group
- The bandwidth properties of each priority group and priority

When you configure hierarchical scheduling on a port, any priority that is not part of an explicitly configured priority group is assigned to the automatically generated default priority group and receives no bandwidth. The default priority group is transparent. It does not appear in the configuration.

Hierarchical Scheduling Process

Hierarchical scheduling consists of multiple configuration steps that create the priorities and the priority groups, schedule their resources, and assign them to interfaces. The steps below correspond to the six blocks in the packet flow diagram shown in [Figure 6 on page 228](#):

1. Packet classification:

- Configure classification of incoming traffic into forwarding classes (priorities). This consists of either using the default classifiers or configuring classifiers to map code points and loss priorities to the forwarding classes.
- Apply the classifiers to ingress interfaces or use the default classifiers. Applying a classifier to an interface groups incoming traffic on the interface into forwarding classes and loss priorities, by applying the classifier code point mapping to the incoming traffic.

2. Configure the output queues for the forwarding classes (priorities). This consists of either using the default forwarding classes and forwarding-class-to-queue mapping, or creating your own forwarding classes and mapping them to output queues.

3. Allocate resources to the forwarding classes:

- Define resources for the priorities. This consists of configuring schedulers to set minimum guaranteed bandwidth, maximum bandwidth, drop profiles for Weighted Random Early Detection (WRED), and bandwidth priority to apply to a forwarding class. Extra bandwidth is shared among queues in proportion to the minimum guaranteed bandwidth (transmit rate) of each queue.
- Map resources to priorities. This consists of mapping forwarding classes to schedulers, using a scheduler map.

4. Configure priority groups. This consists of mapping forwarding classes (priorities) to forwarding class sets (priority groups) to define the priorities that belong to each priority group.

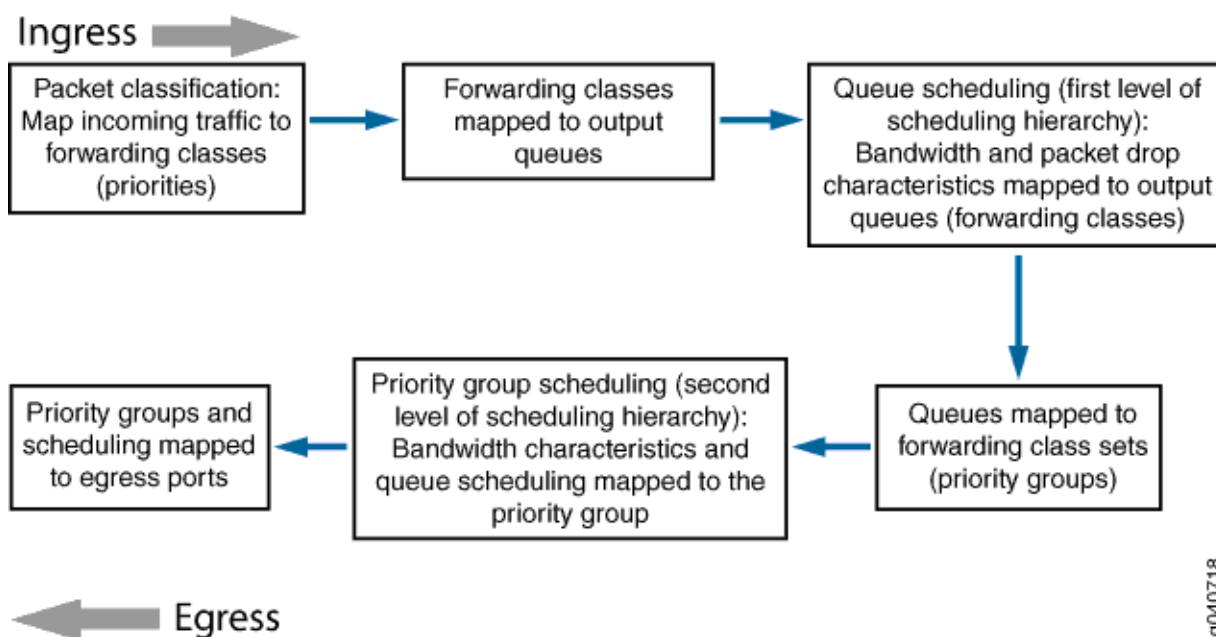
5. Define resources for the priority groups. This consists of configuring traffic control profiles to set minimum guaranteed bandwidth (*guaranteed-rate*) and maximum bandwidth (*shaping-rate* on switches other than QFX10000 switches, *transmit-rate* on QFX10000 switches) for a priority group. Traffic control profiles also specify a scheduler map, which defines the resources (schedulers) mapped to the priorities in the priority group. Extra port bandwidth is shared among priority groups in proportion to the minimum guaranteed bandwidth of each priority group.

The traffic control profile bandwidth settings determine the port resources available to the priority group. The schedulers specified in the scheduler map determine the amount of priority group resources that each priority receives.

NOTE: QFX10000 switches do not support defining a shaping rate for priority groups. Instead, set the maximum bandwidth for a priority group by defining a transmit rate. See *transmit-rate*.

6. Apply hierarchical scheduling to a port. This consists of attaching one or more priority groups (forwarding class sets) to an interface. For each priority group, you also attach a traffic control profile, which contains the scheduling properties of the priority group and the priorities in the priority group. Different priority groups on the same port can use different traffic control profiles, which provides fine tuned control of scheduling for each queue on each interface.

Figure 6: Hierarchical Scheduling Packet Flow



Strict-High Priority Queues and Hierarchical Scheduling

If you configure a strict-high priority queue, you must observe the following rules:

- You must create a separate forwarding class set (priority group) for the strict-high priority queue.
- Only one forwarding class set can contain strict-high priority queues.
- Strict-high priority queues cannot belong to the same forwarding class set as queues that are not strict-high priority.
- A strict-high priority queue cannot belong to a multidestination forwarding class set.
- We recommend that you always apply a *shaping-rate* (*transmit-rate* on QFX10000 switches) to strict-high priority queues to limit the amount of bandwidth a strict-high priority queue can use. If you do not limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

NOTE: On a QFabric system, if a fabric (fte) interface handles strict-high priority traffic, you must define a separate forwarding class set (priority group) for strict-high priority traffic. Strict-high priority traffic cannot be mixed with traffic of other priorities in a forwarding class set. For example, you might choose to create different forwarding class sets for best effort, lossless, strict-high priority, and multidestination traffic.

Default Hierarchical Scheduling

NOTE: There is no default hierarchical scheduling on QFX10000 switches. QFX10000 switches use port scheduling by default, and you must explicitly configure hierarchical scheduling to enable ETS. Also on QFX10000 switches, changing from port scheduler to ETS or from ETS to port scheduler requires a reboot.

If you do not explicitly configure hierarchical scheduling, the switch uses the default settings:

- The switch automatically creates a default forwarding class set that contains all of the forwarding classes on the switch. The switch assigns 100 percent of the port output bandwidth to the default forwarding class set. The default forwarding class set is transparent. It does not appear in the configuration and is used for Data Center Bridging Capability Exchange protocol (DCBX) advertisement.
- Ingress traffic is classified based on the default classifier settings.

- The forwarding classes (queues) in the default forwarding class set receive bandwidth based on the default scheduler settings.

RELATED DOCUMENTATION

Understanding CoS Packet Flow

Understanding CoS Output Queue Schedulers

Understanding CoS Priority Group Scheduling

[Benefits of Configuring CoS Hierarchical Port Scheduling](#)

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Understanding CoS Classifiers

Understanding Default CoS Scheduling and Classification

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Queue Schedulers

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Example: Configuring Minimum Guaranteed Output Bandwidth

Example: Configuring Maximum Output Bandwidth

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

IN THIS SECTION

- [Requirements | 231](#)
- [Overview | 232](#)
- [Configuration | 238](#)
- [Verification | 252](#)

Hierarchical port scheduling defines the class-of-service (CoS) properties of output queues, which are mapped to forwarding classes. Traffic is classified into forwarding classes based on code point (priority), so mapping queues to forwarding classes also maps queues to priorities). Hierarchical port scheduling enables you to group priorities that require similar CoS treatment into priority groups. You define the port bandwidth resources for a priority group, and you define the amount of the priority group's resources that each priority in the group can use.

Hierarchical port scheduling is the Junos OS implementation of enhanced transmission selection (ETS), as described in IEEE 802.1Qaz. One major benefit of hierarchical port scheduling is greater port bandwidth utilization. If a priority group on a port does not use all of its allocated bandwidth, other priority groups on that port can use that bandwidth. Also, if a priority within a priority group does not use its allocated bandwidth, other priorities within that priority group can use that bandwidth.

Configuring hierarchical scheduling is a multistep procedure that includes:

- Mapping forwarding classes to queues
- Defining forwarding class sets (priority groups)
- Defining behavior aggregate classifiers
- Configuring priority-based flow control (PFC) for lossless priorities (queues)
- Applying classifiers and PFC configuration to ingress interfaces
- Defining drop profiles
- Defining schedulers
- Mapping forwarding classes to schedulers
- Defining traffic control profiles
- Assigning priority groups and traffic control profiles to egress ports

NOTE: OCX Series switches do not support lossless transport and do not support PFC. Although this example includes configuring lossless transport with PFC, the portions of the example that do not pertain to lossless transport still apply to OCX Series switches. (You can configure hierarchical scheduling on OCX Series switches, but you cannot configure lossless transport or lossless forwarding classes.)

This example describes how to configure hierarchical scheduling:

Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

Overview

IN THIS SECTION

- [Topology | 233](#)

Keep the following considerations in mind when you plan the port bandwidth allocation for priority groups and for individual priorities:

- How much traffic and what types of traffic you expect to traverse the system.
- How you want to divide different types of traffic into priorities (forwarding classes) to apply different CoS treatment to different types of traffic. Dividing traffic into priorities includes:
 - Mapping the code points of ingress traffic to forwarding classes using behavior aggregate (BA) classifiers. This classifies incoming traffic into the appropriate forwarding class based on code point.
 - Mapping forwarding classes to output queues. This defines the output queue for each type of traffic.
 - Attaching the BA classifier to the desired ingress interfaces so that incoming traffic maps to the desired forwarding classes and queues.
- How you want to organize priorities into priority groups (forwarding class sets).

Traffic that requires similar treatment usually belongs in the same priority group. To do this, place forwarding classes that require similar bandwidth, loss, and other characteristics in the same forwarding class set. For example, you can map all types of best-effort traffic forwarding classes into one forwarding class set.

- How much of the port bandwidth you want to allocate to each priority group and to each of the priorities in each priority group. The following considerations apply to bandwidth allocation:
 - Estimate how much traffic you expect in each forwarding class, and how much traffic you expect in each forwarding class set (the amount of traffic you expect in a forwarding class set is the aggregate amount of traffic in the forwarding classes that belong to the forwarding class set).

- The combined minimum guaranteed bandwidth of the priorities (forwarding classes) in a priority group should not exceed the minimum guaranteed bandwidth of the priority group (forwarding class set). The transmit rate scheduler parameter defines the minimum guaranteed bandwidth for forwarding classes. Scheduler maps associate schedulers with forwarding classes.
- The combined minimum guaranteed bandwidth of the priority groups (forwarding class sets) on a port should not exceed the port's total bandwidth. The guaranteed rate parameter in the traffic control profile defines the minimum bandwidth for a forwarding class set. Associating a scheduler map with a traffic control profile sets the scheduling for the individual forwarding classes in the forwarding class set.

This example creates hierarchical port scheduling by defining priority groups for best effort, guaranteed delivery, and high-performance computing (HPC) traffic. Each priority group includes priorities that need to receive similar CoS treatment. Each priority group and each priority within each priority group receive the CoS resources needed to service their flows. Lossless priorities use PFC to prevent packet loss when the network experiences congestion.

Topology

Table 62 on page 233 shows the configuration components for this example.

NOTE: OCX Series switches do not support lossless transport and do not support PFC. If you eliminate the configuration elements for the default lossless fcoe and no-loss forwarding classes (including classifier, forwarding class set, scheduler, and traffic control profile configuration for those forwarding classes) and for PFC, this example works for OCX Series switches. However, because the default fcoe and no-loss forwarding classes do not carry traffic on OCX Series switches, you can apply the bandwidth allocated to those forwarding classes to other forwarding classes. By default, the active forwarding classes (best-effort, network-control, and mcast) share the unused bandwidth assigned to the fcoe and no-loss forwarding classes.

Table 62: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology

Property	Settings
Hardware	QFX3500 switch

Table 62: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology (*Continued*)

Property	Settings
Mapping of forwarding classes (priorities) to queues	<p>best-effort to queue 0</p> <p>be2 to queue 1</p> <p>fcoe (Fibre Channel over Ethernet) to queue 3</p> <p>no-loss to queue 4</p> <p>hpc (high-performance computing) to queue 5</p> <p>network-control to queue 7</p> <p>NOTE: On switches that do not support the ELS CLI, if you are using Junos OS Release 12.2 or later, use the default forwarding-class-to-queue mapping for the lossless fcoe and no-loss forwarding classes. If you explicitly configure the default lossless forwarding classes, the traffic mapped to those forwarding classes is treated as lossy (best-effort) traffic and does <i>not</i> receive lossless treatment.</p> <p>On switches that do not support the ELS CLI, in Junos OS Release 12.3 and later, you can include the <i>no-loss</i> packet drop attribute in the explicit forwarding class configuration to configure a lossless forwarding class.</p>
Forwarding class sets (priority groups)	<p>best-effort-pg: contains forwarding classes best-effort, be2, and network control</p> <p>guar-delivery-pg: contains forwarding classes fcoe and no-loss</p> <p>hpc-pg: contains forwarding class hpc</p>
Behavior aggregate classifier (maps forwarding classes and loss priorities to incoming packets by IEEE 802.1 code point)	<p>Name—hsclassifier1</p> <p>Code point mapping:</p> <ul style="list-style-type: none"> • 000 to forwarding class best-effort and loss priority low • 001 to forwarding class be2 and loss priority high • 011 to forwarding class fcoe and loss priority low • 100 to forwarding class no-loss and loss priority low • 101 to forwarding class hpc and loss priority low • 110 to forwarding class network-control and loss priority low

Table 62: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology (*Continued*)

Property	Settings
PFC	<p>Congestion notification profile name—gd-cnp</p> <p>PFC enabled on code points: 011 (fcoe priority), 010 (no-loss priority)</p>
Drop profiles	<p>dp-be-low: drop start point 25, drop end point 50, maximum drop rate 80</p> <p>dp-be-high: drop start point 10, drop end point 40, maximum drop rate 100</p> <p>dp-hpc: drop start point 75, drop end point 90, maximum drop rate 75</p> <p>dp-nc: drop start point 80, drop end point 100, maximum drop rate 100</p>
Queue schedulers	<p>be-sched: minimum bandwidth 3g, maximum bandwidth 100%, priority low, drop profiles dp-be-low and dp-be-high</p> <p>fcoe-sched: minimum bandwidth 2.5g, maximum bandwidth 100%, priority low</p> <p>hpc-sched: minimum bandwidth 2g, maximum bandwidth 100%, priority low, drop profile dp-hpc</p> <p>nc-sched: minimum bandwidth 500m, maximum bandwidth 100%, priority low, drop profile dp-nc</p> <p>nl-sched: minimum bandwidth 2g, maximum bandwidth 100%, priority low</p>
Forwarding class-to-scheduler mapping	<p>Scheduler map be-map:</p> <p>Forwarding class best-effort, scheduler be-sched</p> <p>Forwarding class be2, scheduler be-sched</p> <p>Forwarding class network-control, scheduler nc-sched</p> <p>Scheduler map gd-map:</p> <p>Forwarding class fcoe, scheduler fcoe-sched</p> <p>Forwarding class no-loss, scheduler nl-sched</p> <p>Scheduler map hpc-map:</p> <p>Forwarding class hpc, scheduler hpc-sched</p>

Table 62: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology *(Continued)*

Property	Settings
Traffic control profiles	<p>be-tcp: scheduler map be-map, minimum bandwidth 3.5g, maximum bandwidth 100%</p> <p>gd-tcp: scheduler map gd-map, minimum bandwidth 4.5g, maximum bandwidth 100%</p> <p>hpc-tcp: scheduler map hpc-map, minimum bandwidth 2g, maximum bandwidth 100%</p>
Interfaces	<p>This example configures hierarchical port scheduling on interfaces xe-0/0/20 and xe-0/0/21. Because traffic is bidirectional, you apply the ingress and egress configuration components to both interfaces:</p> <ul style="list-style-type: none"> • Classifier Name—hsclassifier1 • Forwarding class sets—best-effort-pg, guar-deliver-pg, hpc-pg • Congestion notification profile—gd-cnp

[Figure 7 on page 237](#) shows a block diagram of the configuration components and the configuration flow of the CLI statements used in the example. You can perform the configuration steps in a different sequence if you want.

Figure 7: Hierarchical Port Scheduling Components Block Diagram

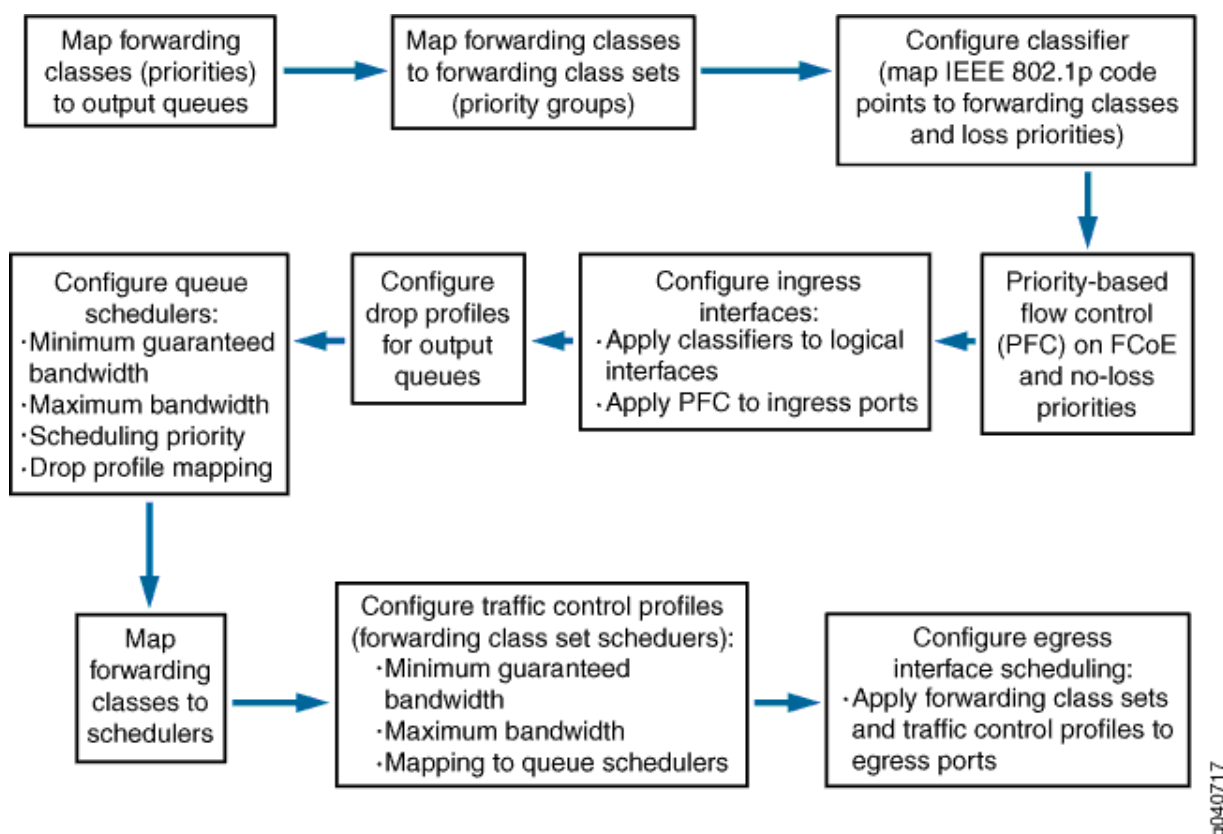
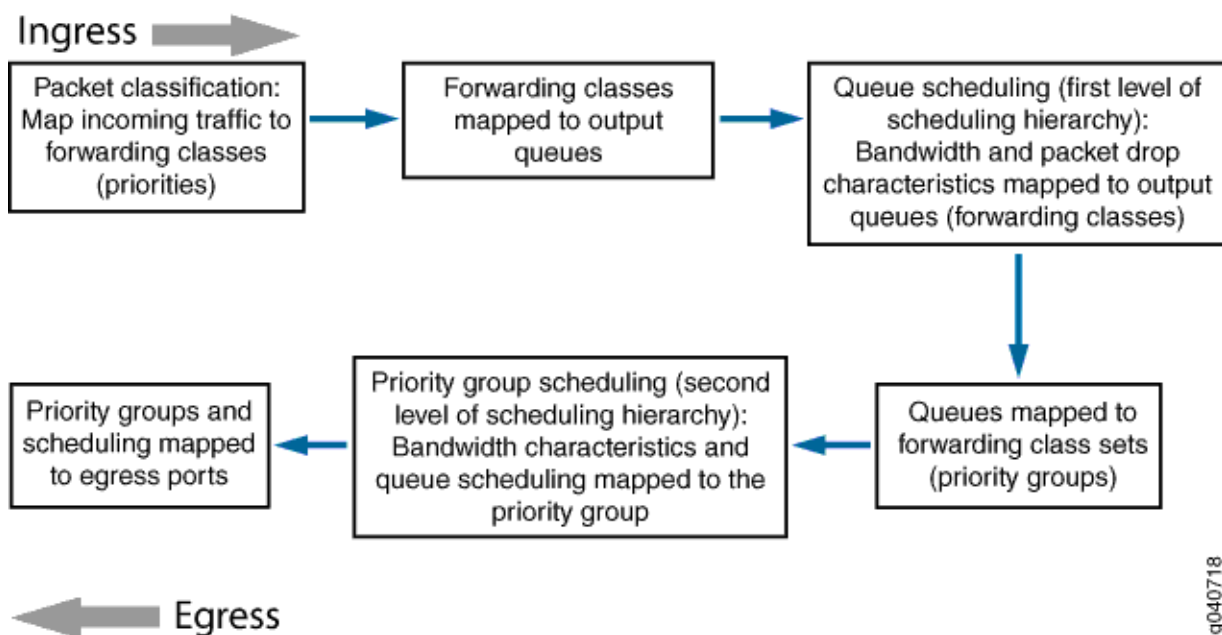


Figure 8 on page 238 shows a block diagram of the hierarchical scheduling packet flow from ingress to egress.

Figure 8: Hierarchical Port Scheduling Packet Flow Block Diagram



Configuration

IN THIS SECTION

- [CLI Quick Configuration | 238](#)
- [Procedure | 242](#)
- [Results | 248](#)

CLI Quick Configuration

To quickly configure hierarchical port scheduling on systems that support lossless transport, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit class-of-service] hierarchy level:

```
[edit class-of-service]
set forwarding-classes class best-effort queue-num 0
set forwarding-classes class be2 queue-num 1
set forwarding-classes class hpc queue-num 5
```

```

set forwarding-classes class network-control queue-num 7
set forwarding-class-sets best-effort-pg class best-effort
set forwarding-class-sets best-effort-pg class be2
set forwarding-class-sets best-effort-pg class network-control
set forwarding-class-sets guar-delivery-pg class fcoe
set forwarding-class-sets guar-delivery-pg class no-loss
set forwarding-class-sets hpc-pg class hpc
set classifiers ieee-802.1 hsclassifier1 forwarding-class best-effort loss-priority low code-
points 000
set classifiers ieee-802.1 hsclassifier1 forwarding-class be2 loss-priority high code-points 001
set classifiers ieee-802.1 hsclassifier1 forwarding-class fcoe loss-priority low code-points
011
set classifiers ieee-802.1 hsclassifier1 forwarding-class no-loss loss-priority low code-points
100
set classifiers ieee-802.1 hsclassifier1 forwarding-class hpc loss-priority low code-points 101
set classifiers ieee-802.1 hsclassifier1 forwarding-class network-control loss-priority low code-
points 110
set congestion-notification-profile gd-cnp input ieee-802.1 code-point 011 pfc
set congestion-notification-profile gd-cnp input ieee-802.1 code-point 100 pfc
set interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 hsclassifier1
set interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 hsclassifier1
set interfaces xe-0/0/20 congestion-notification-profile gd-cnp
set interfaces xe-0/0/21 congestion-notification-profile gd-cnp
set drop-profiles dp-be-low interpolate fill-level 25 fill-level 50 drop-probability 0 drop-
probability 80
set drop-profiles dp-be-high interpolate fill-level 10 fill-level 40 drop-probability 0 drop-
probability 100
set drop-profiles dp-nc interpolate fill-level 80 fill-level 100 drop-probability 0 drop-
probability 100
set drop-profiles dp-hpc interpolate fill-level 75 fill-level 90 drop-probability 0 drop-
probability 75
set schedulers be-sched priority low transmit-rate 3g
set schedulers be-sched shaping-rate percent 100
set schedulers be-sched drop-profile-map loss-priority low protocol any drop-profile dp-be-low
set schedulers be-sched drop-profile-map loss-priority high protocol any drop-profile dp-be-high
set schedulers fcoe-sched priority low transmit-rate 2500m
set schedulers fcoe-sched shaping-rate percent 100
set schedulers hpc-sched priority low transmit-rate 2g
set schedulers hpc-sched shaping-rate percent 100
set schedulers hpc-sched drop-profile-map loss-priority low protocol any drop-profile dp-hpc
set schedulers nc-sched priority low transmit-rate 500m
set schedulers nc-sched shaping-rate percent 100
set schedulers nc-sched drop-profile-map loss-priority low protocol any drop-profile dp-nc

```

```

set schedulers nl-sched priority low transmit-rate 2g
set schedulers nl-sched shaping-rate percent 100
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set scheduler-maps be-map forwarding-class be2 scheduler be-sched
set scheduler-maps be-map forwarding-class network-control scheduler nc-sched
set scheduler-maps gd-map forwarding-class fcoe scheduler fcoe-sched
set scheduler-maps gd-map forwarding-class no-loss scheduler nl-sched
set scheduler-maps hpc-map forwarding-class hpc scheduler hpc-sched
set traffic-control-profiles be-tcp scheduler-map be-map guaranteed-rate 3500m
set traffic-control-profiles be-tcp shaping-rate percent 100
set traffic-control-profiles gd-tcp scheduler-map gd-map guaranteed-rate 4500m
set traffic-control-profiles gd-tcp shaping-rate percent 100
set traffic-control-profiles hpc-tcp scheduler-map hpc-map guaranteed-rate 2g
set traffic-control-profiles hpc-tcp shaping-rate percent 100
set interfaces xe-0/0/20 forwarding-class-set best-effort-pg output-traffic-control-profile be-
tcp
set interfaces xe-0/0/20 forwarding-class-set guar-delivery-pg output-traffic-control-profile gd-
tcp
set interfaces xe-0/0/20 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp
set interfaces xe-0/0/21 forwarding-class-set best-effort-pg output-traffic-control-profile be-
tcp
set interfaces xe-0/0/21 forwarding-class-set guar-delivery-pg output-traffic-control-profile gd-
tcp
set interfaces xe-0/0/21 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp

```

OCX Series Switches

Because OCX Series switches do not support lossless transport, the following subset of the configuration eliminates the lossless configuration elements and provides hierarchical port scheduling for the best-effort, be2, hpc, and network-control forwarding classes. In addition, on OCX Series switches, you would probably use DSCP classifiers and code points instead of IEEE classifiers and code points. To quickly configure hierarchical port scheduling on an OCX Series switch, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit class-of-service] hierarchy level:

```

[edit class-of-service]
set forwarding-classes class best-effort queue-num 0
set forwarding-classes class be2 queue-num 1
set forwarding-classes class hpc queue-num 5
set forwarding-classes class network-control queue-num 7
set forwarding-class-sets best-effort-pg class best-effort

```

```

set forwarding-class-sets best-effort-pg class be2
set forwarding-class-sets best-effort-pg class network-control

set forwarding-class-sets hpc-pg class hpc
set classifiers ieee-802.1 hsclassifier1 forwarding-class best-effort loss-priority low code-points 000
set classifiers ieee-802.1 hsclassifier1 forwarding-class be2 loss-priority high code-points 001

set classifiers ieee-802.1 hsclassifier1 forwarding-class hpc loss-priority low code-points 101
set classifiers ieee-802.1 hsclassifier1 forwarding-class network-control loss-priority low code-points 110

set interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 hsclassifier1
set interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 hsclassifier1
set drop-profiles dp-be-low interpolate fill-level 25 fill-level 50 drop-probability 0 drop-probability 80
set drop-profiles dp-be-high interpolate fill-level 10 fill-level 40 drop-probability 0 drop-probability 100
set drop-profiles dp-nc interpolate fill-level 80 fill-level 100 drop-probability 0 drop-probability 100
set drop-profiles dp-hpc interpolate fill-level 75 fill-level 90 drop-probability 0 drop-probability 75
set schedulers be-sched priority low transmit-rate 3g
set schedulers be-sched shaping-rate percent 100
set schedulers be-sched drop-profile-map loss-priority low protocol any drop-profile dp-be-low
set schedulers be-sched drop-profile-map loss-priority high protocol any drop-profile dp-be-high
set schedulers hpc-sched priority low transmit-rate 2g
set schedulers hpc-sched shaping-rate percent 100
set schedulers hpc-sched drop-profile-map loss-priority low protocol any drop-profile dp-hpc
set schedulers nc-sched priority low transmit-rate 500m
set schedulers nc-sched shaping-rate percent 100
set schedulers nc-sched drop-profile-map loss-priority low protocol any drop-profile dp-nc
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set scheduler-maps be-map forwarding-class be2 scheduler be-sched
set scheduler-maps be-map forwarding-class network-control scheduler nc-sched
set scheduler-maps hpc-map forwarding-class hpc scheduler hpc-sched
set traffic-control-profiles be-tcp scheduler-map be-map guaranteed-rate 3500m
set traffic-control-profiles be-tcp shaping-rate percent 100
set traffic-control-profiles hpc-tcp scheduler-map hpc-map guaranteed-rate 2g
set traffic-control-profiles hpc-tcp shaping-rate percent 100
set interfaces xe-0/0/20 forwarding-class-set best-effort-pg output-traffic-control-profile be-tcp
set interfaces xe-0/0/20 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp

```

```
set interfaces xe-0/0/21 forwarding-class-set best-effort-pg output-traffic-control-profile be-
tcp
set interfaces xe-0/0/21 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp
```

Procedure

Step-by-Step Procedure

To perform a step-by-step configuration of the forwarding classes (priorities), forwarding class sets (priority groups), classifiers, queue schedulers, PFC, traffic control profiles, and interfaces to set up hierarchical port scheduling (ETS):

1. Configure the forwarding classes (priorities) and map them to unicast output queues (do not explicitly map the `fcoe` and `no-loss` forwarding classes to output queues; use the default configuration):

```
[edit class-of-service]
user@switch# set forwarding-classes class best-effort queue-num 0
user@switch# set forwarding-classes class be2 queue-num 1
user@switch# set forwarding-classes class hpc queue-num 5
user@switch# set forwarding-classes class network-control queue-num 7
```

2. Configure forwarding class sets (priority groups) to group forwarding classes (priorities) that require similar CoS treatment:

```
[edit class-of-service]
user@switch# set forwarding-class-sets best-effort-pg class best-effort
user@switch# set forwarding-class-sets best-effort-pg class be2
user@switch# set forwarding-class-sets best-effort-pg class network-control
user@switch# set forwarding-class-sets guar-delivery-pg class fcoe
user@switch# set forwarding-class-sets guar-delivery-pg class no-loss
user@switch# set forwarding-class-sets hpc-pg class hpc
```

NOTE: On OCX Series switches, you would not configure the `guar-delivery-pg` forwarding class set for lossless traffic.

3. Configure a classifier to set the loss priority and IEEE 802.1 code points assigned to each forwarding class at the ingress:

```
[edit class-of-service]
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class best-effort loss-
priority low code-points 000
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class be2 loss-priority
high code-points 001
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class fcoe loss-priority
low code-points 011
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class no-loss loss-
priority low code-points 100
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class hpc loss-priority
low code-points 101
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class network-control loss-
priority low code-points 110
```

NOTE: On OCX Series switches, you would not configure the fcoe and no-loss portions of the classifier.

4. Configure a congestion notification profile to enable PFC on the FCoE and no-loss queue IEEE 802.1 code points:

```
[edit class-of-service]
user@switch# set congestion-notification-profile gd-cnp input ieee-802.1 code-point 011 pfc
user@switch# set congestion-notification-profile gd-cnp input ieee-802.1 code-point 100 pfc
```

NOTE: This step does not apply to OCX Series switches, which do not support PFC.

5. Assign the classifier to the interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 hsclassifier1
user@switch# set interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 hsclassifier1
```

6. Apply the PFC configuration to the interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 congestion-notification-profile gd-cnp
user@switch# set interfaces xe-0/0/21 congestion-notification-profile gd-cnp
```

NOTE: This step does not apply to OCX Series switches, which do not support PFC.

7. Configure the drop profile for the best-effort low loss-priority queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-be-low interpolate fill-level 25 fill-level 50 drop-
probability 0 drop-probability 80
```

8. Configure the drop profile for the best-effort high loss-priority queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-be-high interpolate fill-level 10 fill-level 40 drop-
probability 0 drop-probability 100
```

9. Configure the drop profile for the network-control queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-nc interpolate fill-level 80 fill-level 100 drop-
probability 0 drop-probability 100
```

10. Configure the drop profile for the high-performance computing queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-hpc interpolate fill-level 75 fill-level 90 drop-
probability 0 drop-probability 75
```

11. Define the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profiles for the best-effort queue:

```
[edit class-of-service]
user@switch# set schedulers be-sched priority low transmit-rate 3g
user@switch# set schedulers be-sched shaping-rate percent 100
user@switch# set schedulers be-sched drop-profile-map loss-priority low protocol any drop-
profile dp-be-low
user@switch# set schedulers be-sched drop-profile-map loss-priority high protocol any drop-
profile dp-be-high
```

12. Define the minimum guaranteed bandwidth, priority, and maximum bandwidth for the FCoE queue:

```
[edit class-of-service]
user@switch# set schedulers fcoe-sched priority low transmit-rate 2500m
user@switch# set schedulers fcoe-sched shaping-rate percent 100
```

NOTE: This step does not apply to OCX Series switches, which do not support lossless transport.

13. Define the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profile for the high-performance computing queue:

```
[edit class-of-service]
user@switch# set schedulers hpc-sched priority low transmit-rate 2g
user@switch# set schedulers hpc-sched shaping-rate percent 100
user@switch# set schedulers hpc-sched drop-profile-map loss-priority low protocol any drop-
profile dp-hpc
```

14. Define the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profile for the network-control queue:

```
[edit class-of-service]
user@switch# set schedulers nc-sched priority low transmit-rate 500m
user@switch# set schedulers nc-sched shaping-rate percent 100
```



```
user@switch# set schedulers nc-sched drop-profile-map loss-priority low protocol any drop-profile dp-nc
```

15. Define the minimum guaranteed bandwidth, priority, and maximum bandwidth for the no-loss queue:

```
[edit class-of-service]
user@switch# set schedulers nl-sched priority low transmit-rate 2g
user@switch# set schedulers nl-sched shaping-rate percent 100
```

NOTE: This step does not apply to OCX Series switches, which do not support lossless transport.

16. Map the schedulers to the appropriate forwarding classes (queues):

```
[edit class-of-service]
user@switch# set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
user@switch# set scheduler-maps be-map forwarding-class be2 scheduler be-sched
user@switch# set scheduler-maps be-map forwarding-class network-control scheduler nc-sched
user@switch# set scheduler-maps gd-map forwarding-class fcoe scheduler fcoe-sched
user@switch# set scheduler-maps gd-map forwarding-class no-loss scheduler nl-sched
user@switch# set scheduler-maps hpc-map forwarding-class hpc scheduler hpc-sched
```

NOTE: On OCX Series switches, because lossless transport is not supported, you would not configure the `gd-map` scheduler map.

17. Define the traffic control profile for the best-effort priority group (queue scheduler to mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
user@switch# set traffic-control-profiles be-tcp scheduler-map be-map guaranteed-rate 3500m
user@switch# set traffic-control-profiles be-tcp shaping-rate percent 100
```

18. Define the traffic control profile for the guaranteed delivery priority group (queue to scheduler mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
user@switch# set traffic-control-profiles gd-tcp scheduler-map gd-map guaranteed-rate 4500m
user@switch# set traffic-control-profiles gd-tcp shaping-rate percent 100
```

NOTE: This step does not apply to OCX Series switches, which do not support lossless transport.

19. Define the traffic control profile for the high-performance computing priority group (queue to scheduler mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
user@switch# set traffic-control-profiles hpc-tcp scheduler-map hpc-map guaranteed-rate 2g
user@switch# set traffic-control-profiles hpc-tcp shaping-rate percent 100
```

20. Apply the three priority groups (forwarding class sets) and the appropriate traffic control profiles to the egress ports:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 forwarding-class-set best-effort-pg output-traffic-control-profile be-tcp
user@switch# set interfaces xe-0/0/20 forwarding-class-set guar-delivery-pg output-traffic-control-profile gd-tcp
user@switch# set interfaces xe-0/0/20 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp
user@switch# set interfaces xe-0/0/21 forwarding-class-set best-effort-pg output-traffic-control-profile be-tcp
user@switch# set interfaces xe-0/0/21 forwarding-class-set guar-delivery-pg output-traffic-control-profile gd-tcp
user@switch# set interfaces xe-0/0/21 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp
```

NOTE: Because OCX Series switches do not support lossless transport, on OCX Series switches, you would not apply the `guar-deliver-pg` forwarding class set and the `gd-tcp` traffic control profile to interfaces.

Results

Display the results of the configuration (the system shows only the explicitly configured parameters; it does not show default parameters such as the `fcoe` and `no-loss` lossless forwarding classes). On OCX Series switches, you would not see the lossless configuration components in the output:

```
user@switch> show configuration class-of-service
classifiers {
    ieee-802.1 hsclassifier1 {
        forwarding-class best-effort {
            loss-priority low code-points 000;
        }
        forwarding-class be2 {
            loss-priority high code-points 001;
        }
        forwarding-class fcoe {
            loss-priority low code-points 011;
        }
        forwarding-class no-loss {
            loss-priority low code-points 100;
        }
        forwarding-class hpc {
            loss-priority low code-points 101;
        }
        forwarding-class network-control {
            loss-priority low code-points 110;
        }
    }
}
drop-profiles {
    dp-be-low {
        interpolate {
            fill-level [ 25 50 ];
            drop-probability [ 0 80 ];
        }
    }
}
```

```

dp-be-high {
    interpolate {
        fill-level [ 10 40 ];
        drop-probability [ 0 100 ];
    }
}
dp-hpc {
    interpolate {
        fill-level [ 75 90 ];
        drop-probability [ 0 75 ];
    }
}
dp-nc {
    interpolate {
        fill-level [ 80 100 ];
        drop-probability [ 0 100 ];
    }
}
}
forwarding-classes {
    class best-effort queue-num 0;
    class be2 queue-num 1;
    class hpc queue-num 5;
    class network-control queue-num 7;
}
traffic-control-profiles {
    be-tcp {
        scheduler-map be-map;
        shaping-rate percent 100;
        guaranteed-rate 3500000000;
    }
    gd-tcp {
        scheduler-map gd-map;
        shaping-rate percent 100;
        guaranteed-rate 4500000000;
    }
    hpc-tcp {
        scheduler-map hpc-map;
        shaping-rate percent 100;
        guaranteed-rate 2g;
    }
}
forwarding-class-sets {

```

```

    guar-delivery-pg {
        class fcoe;
        class no-loss;
    }
    best-effort-pg {
        class best-effort;
        class be2;
        class network-control;
    }
    hpc-pg {
        class hpc;
    }
}
congestion-notification-profile {
    gd-cnp {
        input {
            ieee-802.1 {
                code-point 011 {
                    pfc;
                }
                code-point 100 {
                    pfc;
                }
            }
        }
    }
}
}
interfaces {
    xe-0/0/20 {
        forwarding-class-set {
            best-effort-pg {
                output-traffic-control-profile be-tcp;
            }
            guar-delivery-pg {
                output-traffic-control-profile gd-tcp;
            }
            hpc-pg {
                output-traffic-control-profile hpc-tcp;
            }
        }
        congestion-notification-profile gd-cnp;
        unit 0 {
            classifiers {

```

```

        ieee-802.1 hsclassifier1;
    }
}
xe-0/0/21 {
    forwarding-class-set {
        best-effort-pg {
            output-traffic-control-profile be-tcp;
        }
        guar-delivery-pg {
            output-traffic-control-profile gd-tcp;
        }
        hpc-pg {
            output-traffic-control-profile hpc-tcp;
        }
    }
    congestion-notification-profile gd-cnp;
    unit 0 {
        classifiers {
            ieee-802.1 hsclassifier1;
        }
    }
}
scheduler-maps {
    be-map {
        forwarding-class best-effort scheduler be-sched;
        forwarding-class network-control scheduler nc-sched;
        forwarding-class be2 scheduler be-sched;
    }
    gd-map {
        forwarding-class fcoe scheduler fcoe-sched;
        forwarding-class no-loss scheduler nl-sched;
    }
    hpc-map {
        forwarding-class hpc scheduler hpc-sched;
    }
}
schedulers {
    be-sched {
        transmit-rate 3g;
        shaping-rate percent 100;
        priority low;
    }
}

```

```

        drop-profile-map loss-priority low protocol any drop-profile dp-be-low;
        drop-profile-map loss-priority high protocol any drop-profile dp-be-high;
    }
    fcoe-sched {
        transmit-rate 2500000000;
        shaping-rate percent 100;
        priority low;
    }
    hpc-sched {
        transmit-rate 2g;
        shaping-rate percent 100;
        priority low;
        drop-profile-map loss-priority low protocol any drop-profile dp-hpc;
    }
    nc-sched {
        transmit-rate 500m;
        shaping-rate percent 100;
        priority low;
        drop-profile-map loss-priority low protocol any drop-profile dp-nc;
    }
    nl-sched {
        transmit-rate 2g;
        shaping-rate percent 100;
        priority low;
    }
}

```

TIP: To quickly configure the interfaces, issue the `load merge` terminal command, and then copy the hierarchy and paste it into the switch terminal window.

Verification

IN THIS SECTION

- [Verifying the Forwarding Classes \(Priorities\) | 253](#)
- [Verifying the Forwarding Class Sets \(Priority Groups\) | 254](#)
- [Verifying the Classifier | 255](#)
- [Verifying Priority-Based Flow Control | 256](#)

- [Verifying the Output Queue Schedulers | 257](#)
- [Verifying the Drop Profiles | 261](#)
- [Verifying the Priority Group Output Schedulers \(Traffic Control Profiles\) | 262](#)
- [Verifying the Interface Configuration | 263](#)

NOTE: The verification output is based on the full example configuration. On OCX Series switches, you do not see lossless configuration components in the output. Comments about lossless configuration components do not apply to OCX Series switches.

To verify that you created the hierarchical port scheduling components and they are operating properly, perform these tasks:

Verifying the Forwarding Classes (Priorities)

Purpose

Verify that you created the forwarding classes and mapped them to the correct queues. (The system shows only the explicitly configured forwarding classes. It does not show default forwarding classes such as fcoe and no-loss.)

Action

List the forwarding classes using the operational mode command `show class-of-service forwarding-class`:

```
user@switch> show class-of-service forwarding-class
```

Forwarding class	ID	Queue	Policing priority	No-Loss
best-effort	0	0	normal	Disabled
be2	1	3	normal	Disabled
hpc	2	4	normal	Disabled
network-control	3	7	normal	Disabled
mcast	8	8	normal	Disabled

Meaning

The `show class-of-service forwarding-class` command lists all of the configured forwarding classes, the internal identification number of each forwarding class, the queues that are mapped to the forwarding classes, the policing priority, and whether the forwarding class is lossless (no-loss packet drop attribute enabled) or lossy forwarding class (no-loss packet drop attribute disabled). The command output shows that:

- Forwarding class `best-effort` maps to queue 0 and is lossy
- Forwarding class `be2` maps to queue 1 and is lossy
- Forwarding class `hpc` maps to queue 5 and is lossy
- Forwarding class `network-control` maps to queue 7 and is lossy

In addition, the command lists the default multicast (multidestination) forwarding class and the default queue to which it is mapped.

Verifying the Forwarding Class Sets (Priority Groups)

Purpose

Verify that you created the priority groups and that the correct priorities (forwarding classes) belong to the appropriate priority group.

Action

List the forwarding class sets using the operational mode command `show class-of-service forwarding-class-set`:

```
user@switch> show class-of-service forwarding-class-set
Forwarding class set: best-effort-pg, Type: normal-type, Forwarding class set index: 19907
  Forwarding class      Index
  best-effort           0
  be2                   1
  network-control       5

Forwarding class set: guar-delivery-pg, Type: normal-type, Forwarding class set index: 43700
  Forwarding class      Index
  fcoe                  2
  no-loss               3
```

```
Forwarding class set: hpc-pg, Type: normal-type, Forwarding class set index: 60758
```

Forwarding class	Index
hpc	4

Meaning

The `show class-of-service forwarding-class-set` command lists all of the configured forwarding class sets (priority groups), the forwarding classes (priorities) that belong to each priority group, and the internal index number of each priority group. The command output shows that:

- The forwarding class set `best-effort-pg` includes the forwarding classes `best-effort`, `be2`, and `network-control`.
- The forwarding class set `guar-delivery-pg` includes the forwarding classes `fcoe` and `no-loss`.
- The forwarding class set `hpc-pg` includes the forwarding class `hpc`.

Verifying the Classifier

Purpose

Verify that the classifier maps forwarding classes to the correct IEEE 802.1p code points and packet loss priorities.

Action

List the classifier configured for hierarchical port scheduling using the operational mode command `show class-of-service classifier name hsclassifier1`:

```
user@switch> show class-of-service classifier name hsclassifier1
Classifier: hsclassifier1, Code point type: ieee-802.1, Index: 43607
  Code point      Forwarding class      Loss priority
  000             best-effort             low
  001             be2                  high
  011             fcoe                  low
  100             no-loss                low
  101             hpc                  low
  110             network-control        low
```

Meaning

The `show class-of-service classifier name hsclassifier1` command lists all of the IEEE 802.1p code points and the loss priorities mapped to all of the forwarding classes in the classifier. The command output shows that the forwarding classes `best-effort`, `be2`, `no-loss`, `fcoe`, `hpc`, and `network-control` have been created and mapped to IEEE 802.1p code points and loss priorities.

Verifying Priority-Based Flow Control

Purpose

Verify that PFC is enabled on the correct priorities for lossless transport.

Action

List the congestion notification profiles using the operational mode command `show class-of-service congestion-notification`:

```
user@switch> show class-of-service congestion-notification
Type: Input, Name: gd-cnp, Index: 51687
Cable Length: 100 m
  Priority    PFC      MRU
  000        Disabled
  001        Disabled
  010        Disabled
  011        Enabled   2500
  100        Enabled   2500
  101        Disabled
  110        Disabled
  111        Disabled
Type: Output
  Priority    Flow-Control-Queues
  000
      0
  001
      1
  010
      2
  011
      3
  100
```

	4
101	
	5
110	
	6
111	
	7

Meaning

The `show class-of-service congestion-notification` command lists all of the congestion notification profiles and the IEEE 802.1p code points with PFC enabled. The command output shows that PFC is enabled for code points 011 (fcoe priority and queue) and 100 (no-loss priority and queue) for the `gd-cnp` congestion notification profile.

The command also shows the default cable length (100 meters), the default maximum receive unit (2500 bytes), and the default mapping of priorities to output queues because this example does not include configuring these options.

Verifying the Output Queue Schedulers

Purpose

Verify that you created the output queue schedulers with the correct bandwidth parameters and priorities, mapped to the correct queues, and mapped to the correct drop profiles.

Action

List the scheduler maps using the operational mode command `show class-of-service scheduler-map`:

```
user@switch> show class-of-service scheduler-map
Scheduler map: be-map, Index: 64023

Scheduler: be-sched, Forwarding class: best-effort, Index: 13005
  Transmit rate: 3000000000 bps, Rate Limit: none, Buffer size: remainder,
  Buffer Limit: none, Priority: low
  Excess Priority: unspecified
  Shaping rate: 100 percent,
  drop-profile-map-set-type: mark
  Drop profiles:
    Loss priority  Protocol    Index    Name
```

Low	any	55387	dp-be-low
Medium high	any	1	<default-drop-profile>
High	any	4369	dp-be-high

Scheduler: be-sched, Forwarding class: be2, Index: 13005

Transmit rate: 3000000000 bps, Rate Limit: none, Buffer size: remainder,

Buffer Limit: none, Priority: low

Excess Priority: unspecified

Shaping rate: 100 percent,

drop-profile-map-set-type: mark

Drop profiles:

Loss priority	Protocol	Index	Name
Low	any	55387	dp-be-low
Medium high	any	1	<default-drop-profile>
High	any	4369	dp-be-high

Scheduler: nc-sched, Forwarding class: network-control, Index: 45740

Transmit rate: 500000000 bps, Rate Limit: none, Buffer size: remainder,

Buffer Limit: none, Priority: low

Excess Priority: unspecified

Shaping rate: 100 percent,

drop-profile-map-set-type: mark

Drop profiles:

Loss priority	Protocol	Index	Name
Low	any	44207	dp-nc
Medium high	any	1	<default-drop-profile>
High	any	1	<default-drop-profile>

Scheduler map: gd-map, Index: 61447

Scheduler: fcoe-sched, Forwarding class: fcoe, Index: 37289

Transmit rate: 2500000000 bps, Rate Limit: none, Buffer size: remainder,

Buffer Limit: none, Priority: low

Excess Priority: unspecified

Shaping rate: 100 percent,

drop-profile-map-set-type: mark

Drop profiles:

Loss priority	Protocol	Index	Name
Low	any	44207	<default-drop-profile>
Medium high	any	1	<default-drop-profile>
High	any	1	<default-drop-profile>

Scheduler: nl-sched, Forwarding class: no-loss, Index: 29359

```

Transmit rate: 2000000000 bps, Rate Limit: none, Buffer size: remainder,
Buffer Limit: none, Priority: low
Excess Priority: unspecified
Shaping rate: 100 percent,
drop-profile-map-set-type: mark
Drop profiles:
  Loss priority  Protocol  Index  Name
  Low            any       44207  <default-drop-profile>
  Medium high    any       1      <default-drop-profile>
  High           any       1      <default-drop-profile>

```

Scheduler map: hpc-map, Index: 56941

```

Scheduler: hpc-sched, Forwarding class: hpc, Index: 55900
Transmit rate: 2000000000 bps, Rate Limit: none, Buffer size: remainder,
Buffer Limit: none, Priority: low
Excess Priority: unspecified
Shaping rate: 100 percent,
drop-profile-map-set-type: mark
Drop profiles:
  Loss priority  Protocol  Index  Name
  Low            any       57716  dp-hpc
  Medium high    any       1      <default-drop-profile>
  High           any       1      <default-drop-profile>

```

Meaning

The `show class-of-service scheduler-map` command lists all of the configured scheduler maps. For each scheduler map, the command output includes:

- The name of the scheduler map (scheduler-map field)
- The name of the scheduler (scheduler field)
- The forwarding classes mapped to the scheduler (forwarding-class field)
- The minimum guaranteed queue bandwidth (transmit-rate field)
- The scheduling priority (priority field)
- The maximum bandwidth in the priority group the queue can consume (shaping-rate field)
- The drop profile loss priority (loss priority field) for each drop profile name (name field)

The command output shows that:

- The scheduler map `be-map` was created and has these properties:
 - There are two schedulers, `be-sched` and `nc-sched`.
 - The scheduler `be-sched` has two forwarding classes, `best-effort` and `be2`.
 - Scheduler `be-sched` forwarding classes `best-effort` and `be2` share a minimum guaranteed bandwidth of 3,000,000,000 bps, can consume a maximum of 100 percent of the priority group bandwidth, and use the drop profile `dp-be-low` for low loss-priority traffic, the default drop profile for medium-high loss-priority traffic, and the drop profile `dp-be-high` for high loss-priority traffic.
 - The scheduler `nc-sched` has one forwarding class, `network-control`.
 - The `network-control` forwarding class has a minimum guaranteed bandwidth of 500,000,000 bps, can consume a maximum of 100 percent of the priority group bandwidth, and uses the drop profile `dp-nc` for low loss-priority traffic and the default drop profile for medium-high and high loss priority traffic.
- The scheduler map `gd-map` was created and has these properties:
 - There are two schedulers, `fcoe-sched` and `n1-sched`.
 - The scheduler `fcoe-sched` has one forwarding class, `fcoe`.
 - The `fcoe` forwarding class has a minimum guaranteed bandwidth of 2,500,000,000 bps, and can consume a maximum of 100 percent of the priority group bandwidth.
 - The scheduler `n1-sched` has one forwarding class, `no-loss`.
 - The `no-loss` forwarding class has a minimum guaranteed bandwidth of 2,000,000,000 bps, and can consume a maximum of 100 percent of the priority group bandwidth.
- The scheduler map `hpc-map` was created and has these properties:
 - There is one scheduler, `hpc-sched`.
 - The scheduler `hpc-sched` has one forwarding class, `hpc`.
 - The `hpc` forwarding class has a minimum guaranteed bandwidth of 2,000,000,000 bps, can consume a maximum of 100 percent of the priority group bandwidth, and uses the drop profile `dp-hpc` for low loss-priority traffic and the default drop profile for medium-high and high loss-priority traffic.

Verifying the Drop Profiles

Purpose

Verify that you created the drop profiles dp-be-high, dp-be-low, dp-hpc, and dp-nc with the correct fill levels and drop probabilities.

Action

List the drop profiles using the operational mode command `show configuration class-of-service drop-profiles`:

```
user@switch> show configuration class-of-service drop-profiles
dp-be-low {
    interpolate {
        fill-level [ 25 50 ];
        drop-probability [ 0 80 ];
    }
}
dp-be-high {
    interpolate {
        fill-level [ 10 40 ];
        drop-probability [ 0 100 ];
    }
}
dp-hpc {
    interpolate {
        fill-level [ 75 90 ];
        drop-probability [ 0 75 ];
    }
}
dp-nc {
    interpolate {
        fill-level [ 80 100 ];
        drop-probability [ 0 100 ];
    }
}
```


Meaning

The `show configuration class-of-service drop-profiles` command lists the drop profiles and their properties. The command output shows that there are four drop profiles configured, `dp-be-high`, `dp-be-low`, `dp-hpc`, and `dp-nc`. The output also shows that:

- For `dp-be-low`, the drop start point (the first fill level) is when the queue is 25 percent filled, the drop end point (the second fill level) occurs when the queue is 50 percent filled, and the drop probability at the drop end point is 80 percent.
- For `dp-be-high`, the drop start point (the first fill level) is when the queue is 10 percent filled, the drop end point (the second fill level) occurs when the queue is 40 percent filled, and the drop probability at the drop end point is 100 percent.
- For `dp-hpc`, the drop start point (the first fill level) is when the queue is 75 percent filled, the drop end point (the second fill level) occurs when the queue is 90 percent filled, and the drop probability at the drop end point is 75 percent.
- For `dp-nc`, the drop start point (the first fill level) is when the queue is 80 percent filled, the drop end point (the second fill level) occurs when the queue is 100 percent filled, and the drop probability at the drop end point is 100 percent.

Verifying the Priority Group Output Schedulers (Traffic Control Profiles)

Purpose

Verify that you created the traffic control profiles `be-tcp`, `gd-tcp`, and `hpc-tcp` with the correct bandwidth parameters and scheduler mapping.

Action

List the traffic control profiles using the operational mode command `show class-of-service traffic-control-profile`:

```
user@switch> show class-of-service traffic-control-profile
Traffic control profile: be-tcp, Index: 40535
  Shaping rate: 100 percent
  Scheduler map: be-map
  Guaranteed rate: 3500000000

Traffic control profile: gd-tcp, Index: 37959
  Shaping rate: 100 percent
  Scheduler map: gd-map
```

```
Guaranteed rate: 4500000000
```

```
Traffic control profile: hpc-tcp, Index: 47661
```

```
Shaping rate: 100 percent
```

```
Scheduler map: hpc-map
```

```
Guaranteed rate: 2000000000
```

Meaning

The `show class-of-service traffic-control-profile` command lists all of the configured traffic control profiles. For each traffic control profile, the command output includes:

- The name of the traffic control profile (traffic-control-profile)
- The maximum port bandwidth the priority group can consume (shaping-rate)
- The scheduler map associated with the traffic control profile (scheduler-map)
- The minimum guaranteed priority group port bandwidth (guaranteed-rate)

The command output shows that:

- The traffic control profile `be-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `be-map`, and has a minimum guaranteed bandwidth of 3,500,000,000 bps.
- The traffic control profile `gd-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `gd-map`, and has a minimum guaranteed bandwidth of 4,500,000,000 bps.
- The traffic control profile `hpc-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `hpc-map`, and has a minimum guaranteed bandwidth of 2,000,000,000 bps.

Verifying the Interface Configuration

Purpose

Verify that the classifier, the congestion notification profile, and the forwarding class sets are configured on interfaces `xe-0/0/20` and `xe-0/0/21`.

Action

List the interfaces using the operational mode commands `show configuration class-of-service interfaces xe-0/0/20` and `show configuration class-of-service interfaces xe-0/0/21`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/20
forwarding-class-set {
    best-effort-gp {
        output-traffic-control-profile be-tcp;
    }
    guar-delivery-pg {
        output-traffic-control-profile gd-tcp;
    }
    hpc-pg {
        output-traffic-control-profile hpc-tcp;
    }
}
congestion-notification-profile gd_cnp;
unit 0 {
    classifiers {
        ieee-802.1 hsclassifier1;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/21
forwarding-class-set {
    best-effort-gp {
        output-traffic-control-profile be-tcp;
    }
    guar-delivery-pg {
        output-traffic-control-profile gd-tcp;
    }
    hpc-pg {
        output-traffic-control-profile hpc-tcp;
    }
}
congestion-notification-profile gd_cnp;
unit 0 {
    classifiers {
        ieee-802.1 hsclassifier1;
    }
}
```

```
}
}
```

Meaning

The `show configuration class-of-service interfaces interface-name` command shows that each interface includes the forwarding class sets `best-effort-pg`, `guar-delivery-pg`, and `hpc-pg`, congestion notification profile `gd-cnp`, and the IEEE 802.1p classifier `hsclassifier1`.

RELATED DOCUMENTATION

Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)

[Benefits of Configuring CoS Hierarchical Port Scheduling](#)

Assigning CoS Components to Interfaces

Example: Configuring WRED Drop Profiles

Example: Configuring Drop Profile Maps

Example: Configuring Forwarding Classes

Example: Configuring Forwarding Class Sets

Example: Configuring Queue Schedulers

Example: Configuring Queue Scheduling Priority

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Example: Configuring Minimum Guaranteed Output Bandwidth

Example: Configuring Maximum Output Bandwidth

Configuring CoS PFC (Congestion Notification Profiles)

[Overview of CoS Changes Introduced in Junos OS Release 12.2 | 67](#)

Understanding CoS Hierarchical Port Scheduling (ETS)

Understanding CoS Scheduling Behavior and Configuration Considerations

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth

IN THIS SECTION

- [Guaranteeing Bandwidth Using Hierarchical Scheduling](#) | 266
- [Priority Group Guaranteed Rate \(Guaranteed Minimum Bandwidth\)](#) | 268
- [Queue Transmit Rate \(Guaranteed Minimum Bandwidth\)](#) | 268

You can set a guaranteed minimum bandwidth for individual forwarding classes (queues) and for groups of forwarding classes called *forwarding class sets* (priority groups). Setting a minimum guaranteed bandwidth ensures that priority groups and queues receive the bandwidth required to support the expected traffic.

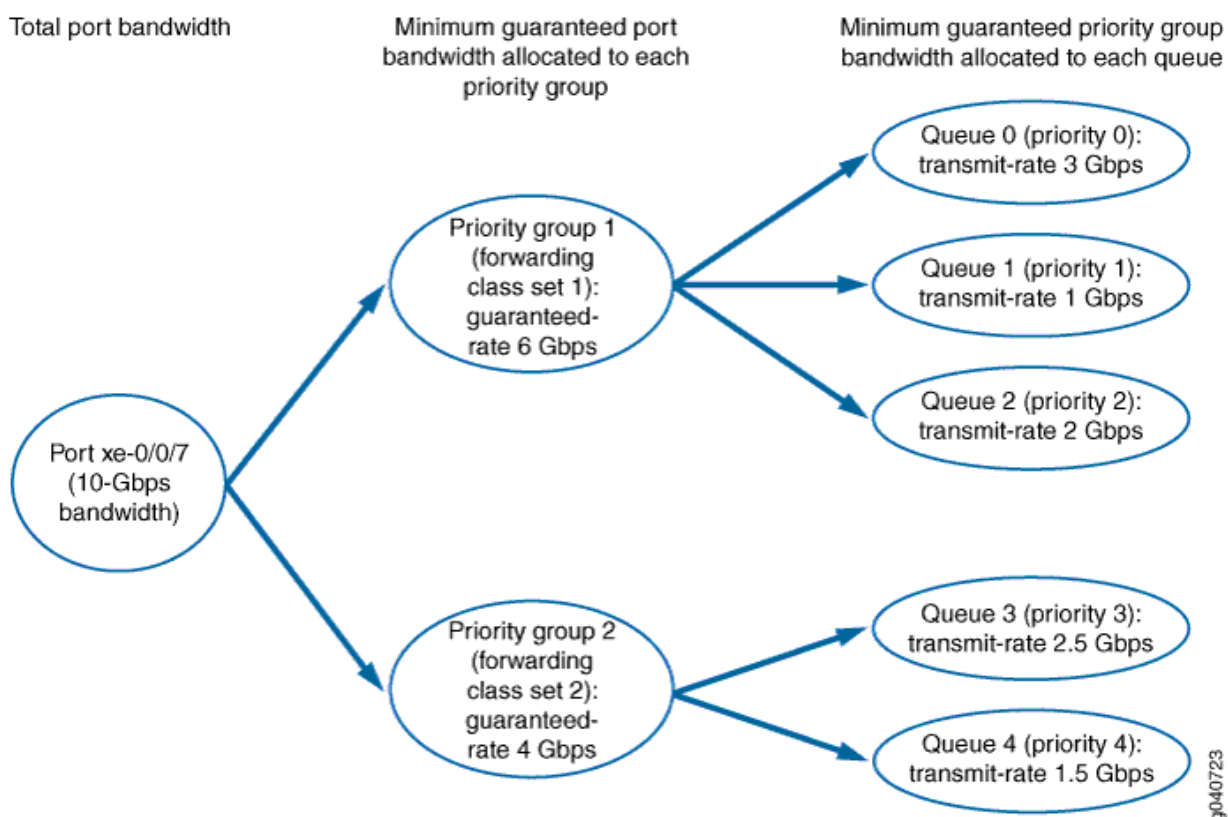
Guaranteeing Bandwidth Using Hierarchical Scheduling

The *guaranteed-rate* value for the priority group (configured in a traffic control profile) defines the minimum amount of bandwidth allocated to a forwarding class set on a port, whereas the *transmit-rate* value of the queue (configured in a scheduler) defines the minimum amount of bandwidth allocated to a particular queue in a priority group. The queue bandwidth is a portion of the priority group bandwidth.

NOTE: You cannot configure a minimum guaranteed bandwidth (transmit rate) for a forwarding class that is mapped to a strict-high priority queue, and you cannot configure a minimum guaranteed bandwidth (guaranteed rate) for a priority group that includes strict-high priority queues.

[Figure 9 on page 267](#) shows how the total port bandwidth is allocated to priority groups (forwarding class sets) based on the guaranteed rate of each priority group. It also shows how the guaranteed bandwidth of each priority group is allocated to the queues in the priority group based on the transmit rate of each queue.

Figure 9: Allocating Guaranteed Bandwidth Using Hierarchical Scheduling



The sum of the priority group guaranteed rates cannot exceed the total port bandwidth. If you configure guaranteed rates whose sum exceeds the port bandwidth, the system sends a syslog message to notify you that the configuration is not valid. However, the system does not perform a commit check. If you commit a configuration in which the sum of the guaranteed rates exceeds the port bandwidth, the hierarchical scheduler behaves unpredictably.

The sum of the queue transmit rates cannot exceed the total guaranteed rate of the priority group to which the queues belong. If you configure transmit rates whose sum exceeds the priority group guaranteed rate, the commit check fails and the system rejects the configuration.

NOTE: You must set both the priority group guaranteed-rate value and the queue transmit-rate value in order to configure the minimum bandwidth for individual queues. If you set the transmit-rate value but do not set the guaranteed-rate value, the configuration fails.

You can set the guaranteed-rate value for a priority group without setting the transmit-rate value for individual queues in the priority group. However, queues that do not have a configured transmit-rate value can become starved for bandwidth if other higher-priority queues need the priority

group's bandwidth. To avoid starving a queue, it is a good practice to configure a transmit-rate value for most queues.

If you configure the guaranteed rate of a priority group as a percentage, configure all of the transmit rates associated with that priority group as percentages. In this case, if any of the transmit rates are configured as absolute values instead of percentages, the configuration is not valid and the system sends a syslog message.

Priority Group Guaranteed Rate (Guaranteed Minimum Bandwidth)

Setting a priority group (forwarding class set) `guaranteed-rate` enables you to reserve a portion of the port bandwidth for the forwarding classes (queues) in that forwarding class set. The minimum bandwidth (`guaranteed-rate`) that you configure for a priority group sets the minimum bandwidth available to all of the forwarding classes in the forwarding class set.

The combined `guaranteed-rate` value of all of the forwarding class sets associated with an interface cannot exceed the amount of bandwidth available on that interface.

You configure the priority group `guaranteed-rate` in the traffic control profile. You cannot apply a traffic control profile that has a guaranteed rate to a priority group that includes a strict-high priority queue.

Queue Transmit Rate (Guaranteed Minimum Bandwidth)

Setting a queue (forwarding class) `transmit-rate` enables you to reserve a portion of the priority group bandwidth for the individual queue. For example, a queue that handles Fibre Channel over Ethernet (FCoE) traffic might require a minimum rate of 4 Gbps to ensure the *class of service* that storage area network (SAN) traffic requires.

The priority group `guaranteed-rate` sets the aggregate minimum amount of bandwidth available to the queues that belong to the priority group. The cumulative total minimum bandwidth the queues consume cannot exceed the minimum bandwidth allocated to the priority group to which they belong. (The combined transmit rates of the queues in a priority group cannot exceed the priority group's guaranteed rate.)

You must configure the `guaranteed-rate` value of the priority group in order to set a `transmit-rate` value for individual queues that belong to the priority group. The reason is that if there is no guaranteed bandwidth for a priority group, there is no way to guarantee bandwidth for queues in that priority group.

You configure the queue `transmit-rate` in the scheduler configuration. You cannot configure a transmit rate for a strict-high priority queue.

RELATED DOCUMENTATION

Understanding CoS Output Queue Schedulers

Understanding CoS Traffic Control Profiles

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Queue Schedulers

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Defining CoS Queue Schedulers

Defining CoS Traffic Control Profiles (Priority Group Scheduling)

Understanding CoS Priority Group Shaping and Queue Shaping (Maximum Bandwidth)

IN THIS SECTION

- [Priority Group Shaping | 269](#)
- [Queue Shaping | 270](#)
- [Shaping Maximum Bandwidth Using Hierarchical Scheduling | 270](#)

If the amount of traffic on an interface exceeds the maximum bandwidth available on the interface, it leads to congestion. You can use priority group (forwarding class set) shaping and queue (forwarding class) shaping to manage traffic and avoid congestion.

Configuring a maximum bandwidth sets the most bandwidth a priority group or a queue can use after all of the priority group and queue minimum bandwidth requirements are met, even if more bandwidth is available.

Priority Group Shaping

Priority group shaping enables you to shape the aggregate traffic of a forwarding class set on a port to a maximum rate that is less than the line or port rate. The maximum bandwidth (*shaping-rate*) that you configure for a priority group sets the maximum bandwidth available to all of the forwarding classes (queues) in the forwarding class set.

If a port has more than one priority group and the combined `shaping-rate` value of the priority groups is greater than the amount of port bandwidth available, the bandwidth is shared proportionally among the priority groups.

You configure the priority group `shaping-rate` in the traffic control profile.

Queue Shaping

Queue shaping throttles the rate at which queues transmit packets. For example, using queue shaping, you can rate-limit a strict-high priority queue so that the strict-priority queue does not lock out (or starve) low-priority queues.

NOTE: We recommend that you always apply a shaping rate to strict-high priority queues to prevent them from starving other queues. If you do not apply a shaping rate to limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

Similarly, for any queue, you can configure queue shaping (`shaping-rate`) to set the maximum bandwidth for a particular queue.

The `shaping-rate` value of the priority group sets the aggregate maximum amount of bandwidth available to the queues that belong to the priority group. On a port, the cumulative total bandwidth the queues consume cannot exceed the maximum bandwidth of the priority group to which they belong.

If a priority group has more than one queue, and the combined `shaping-rate` of the queues is greater than the amount of bandwidth available to the priority group, the bandwidth is shared proportionally among the queues.

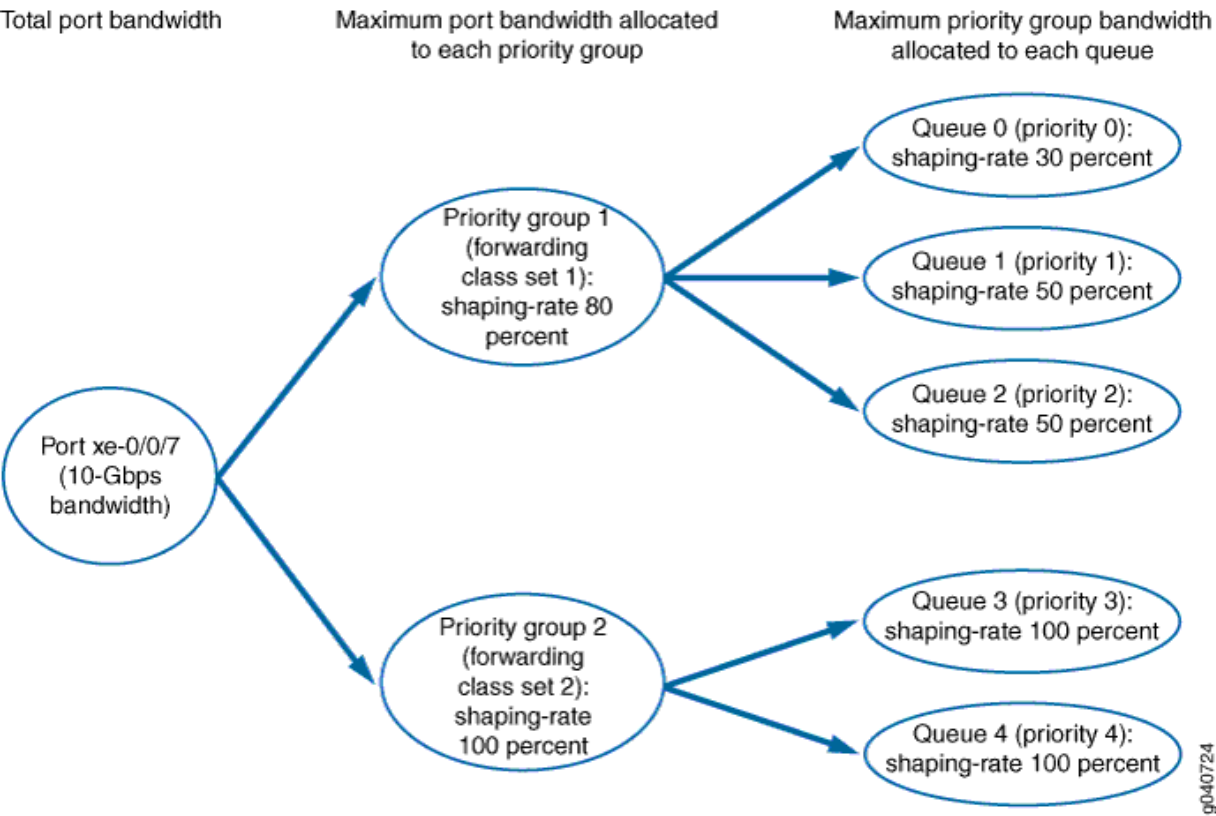
You configure the queue `shaping-rate` in the scheduler configuration, and you set the `shaping-rate` for priority groups in the traffic control profile configuration.

Shaping Maximum Bandwidth Using Hierarchical Scheduling

Priority group shaping defines the maximum bandwidth allocated to a forwarding class set on a port, whereas queue shaping defines a limit on maximum bandwidth usage per queue. The queue bandwidth is a portion of the priority group bandwidth.

[Figure 10 on page 271](#) shows how the port bandwidth is allocated to priority groups (forwarding class sets) based on the shaping rate of each priority group, and how the bandwidth of each priority group is allocated to the queues in the priority group based on the shaping rate of each queue.

Figure 10: Setting Maximum Bandwidth Using Hierarchical Scheduling



RELATED DOCUMENTATION

- [Understanding CoS Output Queue Schedulers](#)
- [Understanding CoS Traffic Control Profiles](#)
- [Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#)
- [Example: Configuring Queue Schedulers](#)
- [Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\)](#)
- [Defining CoS Queue Schedulers](#)
- [Defining CoS Traffic Control Profiles \(Priority Group Scheduling\)](#)

Example: Configuring Minimum Guaranteed Output Bandwidth

IN THIS SECTION

- [Requirements | 273](#)
- [Overview | 274](#)
- [Verification | 276](#)

Scheduling the minimum guaranteed output bandwidth for a queue (forwarding class) requires configuring both tiers of the two-tier hierarchical scheduler. One tier is scheduling the resources for the individual queue. The other tier is scheduling the resources for the priority group (forwarding class set) to which the queue belongs. You set a minimum guaranteed bandwidth to ensure that priority groups and queues receive the bandwidth required to support the expected traffic.

Configuring Guaranteed Minimum Bandwidth

CLI Quick Configuration

To quickly configure the minimum guaranteed bandwidth for a priority group and a queue, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level:

```
[edit class-of-service]
set schedulers be-sched transmit-rate 2g
set traffic-control-profiles be-tcp guaranteed-rate 4g
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set traffic-control-profiles be-tcp scheduler-map be-map
set forwarding-class-sets be-pg class best-effort
set interfaces xe-0/0/7 forwarding-class-set be-pg output-traffic-control-profile be-tcp
```

Step-by-Step Procedure

To configure the minimum guaranteed bandwidth hierarchical scheduling for a queue and a priority group:

1. Configure the minimum guaranteed queue bandwidth of 2 Gbps for scheduler `be-sched`:

```
[edit class-of-service schedulers]
user@switch# set be-sched transmit-rate 2g
```

2. Configure the minimum guaranteed priority group bandwidth of 4 Gbps for traffic control profile `be-tcp`:

```
[edit class-of-service traffic-control-profiles]
user@switch# set be-tcp guaranteed-rate 4g
```

3. Associate the scheduler `be-sched` with the best-effort queue in the scheduler map `be-map`:

```
[edit class-of-service scheduler-maps]
user@switch# set be-map forwarding-class best-effort scheduler be-sched
```

4. Associate the scheduler map with the traffic control profile:

```
[edit class-of-service traffic-control-profiles]
user@switch# set be-tcp scheduler-map be-map
```

5. Assign the best-effort queue to the priority group `be-pg`:

```
[edit class-of-service forwarding-class-sets]
user@switch# set be-pg class best-effort
```

6. Apply the configuration to interface `xe-0/0/7`:

```
[edit class-of-service interfaces]
user@switch# set xe-0/0/7 forwarding-class-set be-pg output-traffic-control-profile be-tcp
```

Requirements

This example uses the following hardware and software components:

- A Juniper Networks QFX3500 Switch

- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

Overview

The priority group minimum guaranteed bandwidth defines the minimum total amount of bandwidth available for all of the queues in the priority group to meet their minimum bandwidth requirements.

The `transmit-rate` setting in the scheduler configuration determines the minimum guaranteed bandwidth for an individual queue. The transmit rate also determines the amount of excess (extra) priority group bandwidth that the queue can share. Extra priority group bandwidth is allocated among the queues in the priority group in proportion to the transmit rate of each queue.

The `guaranteed-rate` setting in the traffic control profile configuration determines the minimum guaranteed bandwidth for a priority group. The guaranteed rate also determines the amount of excess (extra) port bandwidth that the priority group can share. Extra port bandwidth is allocated among the priority groups on a port in proportion to the guaranteed rate of each priority group.

NOTE: You must configure both the `transmit-rate` value for the queue and the `guaranteed-rate` value for the priority group to set a valid minimum bandwidth guarantee for a queue. (If the priority group does not have a guaranteed minimum bandwidth, there is no guaranteed bandwidth pool from which the queue can take its guaranteed minimum bandwidth.)

The sum of the queue transmit rates in a priority group should not exceed the guaranteed rate for the priority group. (You cannot guarantee a minimum bandwidth for the queues that is greater than the minimum bandwidth guaranteed for the entire set of queues.)

NOTE: When you configure bandwidth for a queue or a priority group, the switch considers only the data as the configured bandwidth. The switch does not account for the bandwidth consumed by the preamble and the interframe gap (IFG). Therefore, when you calculate and configure the bandwidth requirements for a queue or for a priority group, consider the preamble and the IFG as well as the data in the calculations.

NOTE: You cannot configure minimum guaranteed bandwidth on strict-high priority queues or on a priority group that contains strict-high priority queues.

This example describes how to:

- Configure a transmit rate (minimum guaranteed queue bandwidth) of 2 Gbps for queues in a scheduler named `be-sched`.

- Configure a guaranteed rate (minimum guaranteed priority group bandwidth) of 4 Gbps for a priority group in a traffic control profile named `be-tcp`.
- Assign the scheduler to a queue named `best-effort` by using a scheduler map named `be-map`.
- Associate the scheduler map `be-map` with the traffic control profile `be-tcp`.
- Assign the queue `best-effort` to a priority group named `be-pg`.
- Assign the priority group and the minimum guaranteed bandwidth scheduling to the egress interface `xe-0/0/7`.

Table 63 on page 275 shows the configuration components for this example:

Table 63: Components of the Minimum Guaranteed Output Bandwidth Configuration Example

Component	Settings
Hardware	QFX3500 switch
Minimum guaranteed queue bandwidth	Transmit rate: 2g
Minimum guaranteed priority group bandwidth	Guaranteed rate: 4g
Scheduler	be-sched
Scheduler map	be-map
Traffic control profile	be-tcp
Forwarding class set (priority group)	be-pg
Queue (forwarding class)	best-effort
Egress interface	xe-0/0/7

Verification

IN THIS SECTION

- [Verifying the Minimum Guaranteed Queue Bandwidth | 276](#)
- [Verifying the Priority Group Minimum Guaranteed Bandwidth and Scheduler Map Association | 276](#)
- [Verifying the Scheduler Map Configuration | 277](#)
- [Verifying Queue \(Forwarding Class\) Membership in the Priority Group | 277](#)
- [Verifying the Egress Interface Configuration | 278](#)

To verify the minimum guaranteed output bandwidth configuration, perform these tasks:

Verifying the Minimum Guaranteed Queue Bandwidth

Purpose

Verify that you configured the minimum guaranteed queue bandwidth as 2g in the scheduler be-sched.

Action

Display the minimum guaranteed bandwidth in the be-sched scheduler configuration using the operational mode command `show configuration class-of-service schedulers be-sched transmit-rate`:

```
user@switch> show configuration class-of-service schedulers be-sched transmit-rate
2g;
```

Verifying the Priority Group Minimum Guaranteed Bandwidth and Scheduler Map Association

Purpose

Verify that the minimum guaranteed priority group bandwidth is 4g and the attached scheduler map is be-map in the traffic control profile be-tcp.

Action

Display the minimum guaranteed bandwidth in the be-tcp traffic control profile configuration using the operational mode command `show configuration class-of-service traffic-control-profiles be-tcp guaranteed-rate`:

```
user@switch> show configuration class-of-service traffic-control-profiles be-tcp guaranteed-rate
4g;
```

Display the scheduler map in the be-tcp traffic control profile configuration using the operational mode command `show configuration class-of-service traffic-control-profiles be-tcp scheduler-map`:

```
user@switch> show configuration class-of-service traffic-control-profiles be-tcp scheduler-map
scheduler-map be-map;
```

Verifying the Scheduler Map Configuration

Purpose

Verify that the scheduler map `be-map` maps the forwarding class `best-effort` to the scheduler `be-sched`.

Action

Display the `be-map` scheduler map configuration using the operational mode command `show configuration class-of-service schedulers maps be-map`:

```
user@switch> show configuration class-of-service scheduler-maps be-map
forwarding-class best-effort scheduler be-sched;
```

Verifying Queue (Forwarding Class) Membership in the Priority Group

Purpose

Verify that the forwarding class set `be-pg` includes the forwarding class `best-effort`.

Action

Display the be-pg forwarding class set configuration using the operational mode command `show configuration class-of-service forwarding-class-sets be-pg`:

```
user@switch> show configuration class-of-service forwarding-class-sets be-pg
class best-effort;
```

Verifying the Egress Interface Configuration

Purpose

Verify that the forwarding class set be-pg and the traffic control profile be-tcp are attached to egress interface xe-0/0/7.

Action

Display the egress interface using the operational mode command `show configuration class-of-service interfaces xe-0/0/7`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/7
forwarding-class-set {
  be-pg {
    output-traffic-control-profile be-tcp;
  }
}
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Queue Schedulers

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Example: Configuring Queue Scheduling Priority

Example: Configuring Forwarding Class Sets

Understanding CoS Traffic Control Profiles

Understanding CoS Hierarchical Port Scheduling (ETS)

Troubleshooting Egress Bandwidth That Exceeds the Configured Minimum Bandwidth

IN THIS SECTION

- [Problem | 279](#)
- [Cause | 279](#)
- [Solution | 280](#)

Problem

Description

The guaranteed minimum bandwidth of a queue (forwarding class) or a priority group (forwarding class set) when measured at the egress port exceeds the guaranteed minimum bandwidth configured for the queue (transmit-rate) or for the priority group (guaranteed-rate).

NOTE: On switches that support enhanced transmission selection (ETS) hierarchical scheduling, the switch allocates guaranteed minimum bandwidth first to a priority group using the guaranteed rate setting in the traffic control profile, and then allocates priority group minimum guaranteed bandwidth to forwarding classes in the priority group using the transmit rate setting in the queue scheduler.

On switches that support direct port scheduling, there is no scheduling hierarchy. The switch allocates port bandwidth to forwarding classes directly, using the transmit rate setting in the queue scheduler.

In this topic, if you are using direct port scheduling on your switch, ignore the references to priority groups and forwarding class sets (priority groups and forwarding class sets are only used for ETS hierarchical port scheduling). For direct port scheduling, only the transmit rate queue scheduler setting can cause the issue described in this topic.

Cause

When you configure bandwidth for a queue or a priority group, the switch accounts for the configured bandwidth as data only. The switch does not include the preamble and the interframe gap (IFG)

associated with frames, so the switch does not account for the bandwidth consumed by the preamble and the IFG in its minimum bandwidth calculations.

The measured egress bandwidth can exceed the configured minimum bandwidth when small packet sizes (64 or 128 bytes) are transmitted because the preamble and the IFG are a larger percentage of the total traffic. For larger packet sizes, the preamble and IFG overhead are a small portion of the total traffic, and the effect on egress bandwidth is minor.

NOTE: For ETS, the sum of the queue transmit rates in a priority group should not exceed the guaranteed rate for the priority group. (You cannot guarantee a minimum bandwidth for the queues that is greater than the minimum bandwidth guaranteed for the entire set of queues.) For port scheduling, the sum of the queue transmit rates should not exceed the port bandwidth.

Solution

When you calculate the bandwidth requirements for queues and priority groups on which you expect a significant amount of traffic with small packet sizes, consider the transmit rate and the guaranteed rate as the minimum bandwidth for the data only. Add sufficient bandwidth to your calculations to account for the preamble and IFG so that the port bandwidth is sufficient to handle the combined minimum data rate and the preamble and IFG.

If the minimum bandwidth measured at the egress port exceeds the amount of bandwidth that you want to allocate to a queue or to a priority group, reduce the transmit rate for that queue and reduce the guaranteed rate of the priority group that contains the queue.

RELATED DOCUMENTATION

[transmit-rate](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth](#)

Example: Configuring Maximum Output Bandwidth

IN THIS SECTION

● [Requirements](#) | 282

- [Overview | 283](#)
- [Verification | 284](#)

Scheduling the maximum output bandwidth for a queue (forwarding class) requires configuring both tiers of the hierarchical scheduler. One tier is scheduling the resources for the individual queue. The other tier is scheduling the resources for the priority group (forwarding class set) to which the queue belongs. You can use priority group and queue shaping to prevent traffic from using more bandwidth than you want the traffic to receive.

Configuring Maximum Bandwidth

CLI Quick Configuration

To quickly configure the maximum bandwidth for a priority group and a queue, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level:

```
[edit class-of-service]
set schedulers be-sched shaping-rate percent 4g
set traffic-control-profiles be-tcp shaping-rate 6g
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set traffic-control-profiles be-tcp scheduler-map be-map
set forwarding-class-sets be-pg class best-effort
set interfaces xe-0/0/7 forwarding-class-set be-pg output-traffic-control-profile be-tcp
```

Step-by-Step Procedure

To configure the maximum bandwidth hierarchical scheduling for a queue and a priority group:

1. Configure the maximum queue bandwidth of 4 Gbps for scheduler be-sched:

```
[edit class-of-service schedulers]
user@switch# set be-sched shaping-rate 4g
```

2. Configure the maximum priority group bandwidth of 6 Gbps for traffic control profile be-tcp:

```
[edit class-of-service traffic-control-profiles]
user@switch# set be-tcp shaping-rate 6g
```

3. Associate the scheduler be-sched with the best-effort queue in the scheduler map be-map:

```
[edit class-of-service scheduler-maps]
user@switch# set be-map forwarding-class best-effort scheduler be-sched
```

4. Associate the scheduler map with the traffic control profile:

```
[edit class-of-service traffic-control-profiles]
user@switch# set be-tcp scheduler-map be-map
```

5. Assign the best-effort queue to the priority group be-pg:

```
[edit class-of-service forwarding-class-sets]
user@switch# set be-pg class best-effort
```

6. Apply the configuration to interface xe-0/0/7:

```
[edit class-of-service interfaces]
user@switch# set xe-0/0/7 forwarding-class-set be-pg output-traffic-control-profile be-tcp
```

Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

Overview

The priority group maximum bandwidth defines the maximum total amount of bandwidth available for all of the queues in the priority group.

The shaping-rate setting in the scheduler configuration determines the maximum bandwidth for an individual queue.

The shaping-rate setting in the traffic control profile configuration determines the maximum bandwidth for a priority group.

NOTE: When you configure bandwidth for a queue or a priority group, the switch considers only the data as the configured bandwidth. The switch does not account for the bandwidth consumed by the preamble and the interframe gap (IFG). Therefore, when you calculate and configure the bandwidth requirements for a queue or for a priority group, consider the preamble and the IFG as well as the data in the calculations.

NOTE: When you set the maximum bandwidth (shaping-rate) for a queue or for a priority group at 100 Kbps or less, the traffic shaping behavior is accurate only within +/- 20 percent of the configured shaping-rate value.

This example describes how to:

- Configure a maximum rate of 4 Gbps for queues in a scheduler named `be-sched`.
- Configure a maximum rate of 6 Gbps for a priority group in a traffic control profile named `be-tcp`.
- Assign the scheduler to a queue named `best-effort` by using a scheduler map named `be-map`.
- Associate the scheduler map `be-map` with the traffic control profile `be-tcp`.
- Assign the queue `best-effort` to a priority group named `be-pg`.
- Assign the priority group and the bandwidth scheduling to the interface `xe-0/0/7`.

Table 64 on page 284 shows the configuration components for this example:

Table 64: Components of the Maximum Output Bandwidth Configuration Example

Component	Settings
Hardware	QFX3500 switch
Maximum queue bandwidth	Shaping rate: 4g
Maximum priority group bandwidth	Shaping rate: 6g
Scheduler	be-sched
Scheduler map	be-map
Traffic control profile	be-tcp
Forwarding class set (priority group)	be-pg
Queue (forwarding class)	best-effort
Egress interface	xe-0/0/7

Verification

IN THIS SECTION

- [Verifying the Maximum Queue Bandwidth | 285](#)
- [Verifying the Priority Group Maximum Bandwidth and Scheduler Map Association | 285](#)
- [Verifying the Scheduler Map Configuration | 286](#)
- [Verifying Queue \(Forwarding Class\) Membership in the Priority Group | 286](#)
- [Verifying the Egress Interface Configuration | 286](#)

To verify the maximum output bandwidth configuration, perform these tasks:

Verifying the Maximum Queue Bandwidth

Purpose

Verify that you configured the maximum queue bandwidth as 4g in the scheduler be-sched.

Action

List the maximum bandwidth in the be-sched scheduler configuration using the operational mode command `show configuration class-of-service schedulers be-sched shaping-rate`:

```
user@switch> show configuration class-of-service schedulers be-sched shaping-rate
4g;
```

Verifying the Priority Group Maximum Bandwidth and Scheduler Map Association

Purpose

Verify that the maximum priority group bandwidth is 6g and the attached scheduler map is be-map in the traffic control profile be-tcp.

Action

List the maximum bandwidth in the be-tcp traffic control profile configuration using the operational mode command `show configuration class-of-service traffic-control-profiles be-tcp shaping-rate`:

```
user@switch> show configuration class-of-service traffic-control-profiles be-tcp shaping-rate
6g;
```

List the scheduler map in the be-tcp traffic control profile configuration using the operational mode command `show configuration class-of-service traffic-control-profiles be-tcp scheduler-map`:

```
user@switch> show configuration class-of-service traffic-control-profiles be-tcp scheduler-map
scheduler-map be-map;
```


Verifying the Scheduler Map Configuration

Purpose

Verify that the scheduler map `be-map` maps the forwarding class `best-effort` to the scheduler `be-sched`.

Action

List the `be-map` scheduler map configuration using the operational mode command `show configuration class-of-service schedulers maps be-map`:

```
user@switch> show configuration class-of-service scheduler-maps be-map
forwarding-class best-effort scheduler be-sched;
```

Verifying Queue (Forwarding Class) Membership in the Priority Group

Purpose

Verify that the forwarding class set `be-pg` includes the forwarding class `best-effort`.

Action

List the `be-pg` forwarding class set configuration using the operational mode command `show configuration class-of-service forwarding-class-sets be-pg`:

```
user@switch> show configuration class-of-service forwarding-class-sets be-pg
class best-effort;
```

Verifying the Egress Interface Configuration

Purpose

Verify that the forwarding class set `be-pg` and the traffic control profile `be-tcp` are attached to egress interface `xe-0/0/7`.

Action

List the egress interface using the operational mode command `show configuration class-of-service interfaces xe-0/0/7`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/7
forwarding-class-set {
    be-pg {
        output-traffic-control-profile be-tcp;
    }
}
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Queue Schedulers

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

Example: Configuring Forwarding Class Sets

Understanding CoS Traffic Control Profiles

Understanding CoS Hierarchical Port Scheduling (ETS)

Troubleshooting Egress Bandwidth That Exceeds the Configured Maximum Bandwidth

IN THIS SECTION

- [Problem | 288](#)
- [Cause | 288](#)
- [Solution | 288](#)

Problem

Description

The maximum bandwidth of a queue when measured at the egress port exceeds the maximum bandwidth rate shaper (shaping-rate statement on QFX5200, QFX5100, EX4600, QFX3500, QFX3600, and OCX1100 switches, and on QFabric systems, and transmit-rate (rate | percentage *percent* exact statement on QFX10000 switches) configured for the queue.

Cause

When you configure bandwidth for a queue (forwarding class) or a priority group (forwarding class set), the switch accounts for the configured bandwidth as data only. The switch does not rate-shape the preamble and the interframe gap (IFG) associated with frames, so the switch does not account for the bandwidth consumed by the preamble and the IFG in its maximum bandwidth calculations.

The measured egress bandwidth can exceed the configured maximum bandwidth when small packet sizes (64 or 128 bytes) are transmitted because the preamble and the IFG are a larger percentage of the total traffic. For larger packet sizes, the preamble and IFG overhead are a small portion of the total traffic, and the effect on egress bandwidth is minor.

Solution

When you calculate the bandwidth requirements for queues on which you expect a significant amount of traffic with small packet sizes, consider the shaping rate as the maximum bandwidth for the data only. Add sufficient bandwidth to your calculations to account for the preamble and IFG so that the port bandwidth is sufficient to handle the combined maximum data rate (shaping rate) and the preamble and IFG.

If the maximum bandwidth measured at the egress port exceeds the amount of bandwidth that you want to allocate to the queue, reduce the shaping rate for that queue.

Troubleshooting Egress Queue Bandwidth Impacted by Congestion

IN THIS SECTION

● [Problem | 289](#)

● [Cause | 289](#)

Problem

Description

Congestion on an egress port causes egress queues to receive less bandwidth than expected. Egress port congestion can impact the amount of bandwidth allocated to queues on the congested port and, in some cases, on ports that are not congested.

Cause

Egress queue congestion can cause the ingress port buffer to fill above a certain threshold and affect the flow to the queues on the egress port. One queue receives its configured bandwidth, but the other queues on the egress port are affected and do not receive their configured share of bandwidth.

Solution

The solution is to configure a drop profile to apply weighted random early detection (WRED) to the queue or queues on the congested ports.

Configure a drop profile on the queue that is receiving its configured bandwidth. This queue is preventing the other queues from receiving their expected bandwidth. The drop profile prevents the queue from affecting the other queues on the port.

To configure a WRED profile using the CLI:

1. Name the drop profile and set the drop start point, drop end point, minimum drop rate, and maximum drop rate for the drop profile:

```
[edit class-of-service]
user@switch# set drop-profile drop-profile-name interpolate fill-level percentage fill-level
percentage drop-probability 0 drop-probability percentage
```

RELATED DOCUMENTATION

| [drop-profile](#)

[Example: Configuring WRED Drop Profiles](#)[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#)[Understanding CoS WRED Drop Profiles](#)

Understanding CoS WRED Drop Profiles

IN THIS SECTION

- [Drop Profile Parameters | 291](#)
- [Defining Drop Profiles on Switches Except QFX10000 | 291](#)
- [Defining Drop Profiles on QFX10000 Switches | 292](#)
- [Default Drop Profile | 293](#)
- [Packet Drop Method | 293](#)
- [Packet Drop Example for Switches Except QFX10000 | 294](#)
- [Drop Profile Maps | 295](#)
- [Congestion Prevention | 295](#)
- [Configuring a WRED Drop Profile and Applying it to an Output Queue | 295](#)
- [Drop Profiles on Explicit Congestion Notification Enabled Queues | 296](#)

When the number of packets queued is greater than the ability of the device to empty an output queue, the queue requires a method for determining which packets to drop to relieve the congestion. Weighted random early detection (WRED) drop profiles define the drop probability of packets of different packet loss probabilities (PLPs) as the output queue fills. During periods of congestion, as the output queue fills, the device drops incoming packets as determined by a drop profile, until the output queue becomes less congested.

Depending on the drop probabilities, a drop profile can drop many packets long before the buffer becomes full, or it can drop only a few packets even if the buffer is almost full.

You configure drop profiles in the drop profile section of the class-of-service (CoS) configuration hierarchy. You apply drop profiles using a drop profile map in queue scheduler configuration. For each queue scheduler, you can configure separate drop profiles for each PLP using the `loss-priority` attribute (low, medium-high, and high). This enables you to treat traffic of different PLPs in different ways during periods of congestion.

NOTE: Do not apply drop profiles to lossless traffic (traffic that belongs to a forwarding class that has the no-loss drop attribute.). Lossless traffic uses priority-based flow control (PFC) to control congestion.

NOTE: You cannot apply drop profiles to multidestination queues on devices that support them.

Drop Profile Parameters

Drop profiles specify two values, which work as pairs:

- **Fill level**—The queue fullness value, which represents a percentage of the memory used to store packets in relation to the total amount of memory allocated to the queue.
- **Drop probability**—The percentage value that corresponds to the likelihood that an individual packet is dropped.

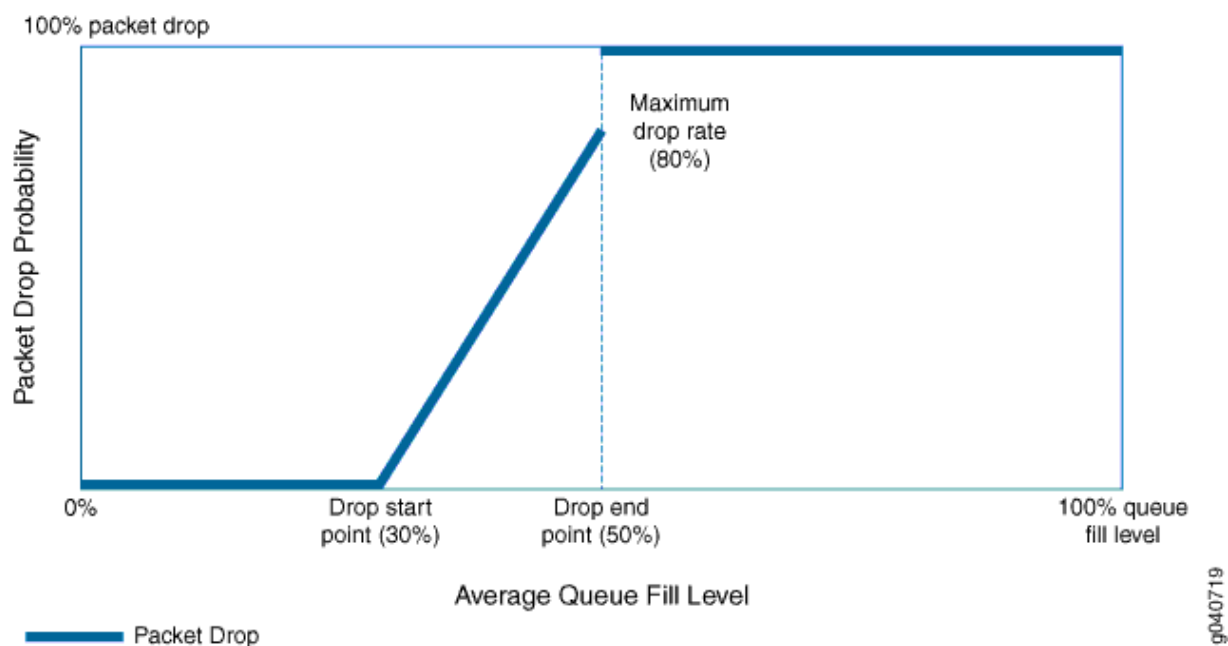
Defining Drop Profiles on Switches Except QFX10000

You set two queue fill levels and two drop probabilities in each drop profile. The first fill level and the first drop probability create one value pair and the second fill level and the second drop probability create a second value pair.

The first fill level value specifies the percentage of queue fullness at which packets begin to drop, known as the drop start point. Until the queue reaches this level of fullness, no packets are dropped. The second fill level value specifies the percentage of queue fullness at which all packets are dropped, known as the drop end point.

The first drop probability value is always 0 (zero). This pairs with the drop start point and specifies that until the queue fullness level reaches the first fill level, no packets drop. When the queue fullness exceeds the drop start point, packets begin to drop until the queue exceeds the second fill level, when all packets drop. The second drop probability value, known as the maximum drop rate, specifies the likelihood of dropping packets when the queue fullness reaches the drop end point. As the queue fills from the drop start point to the drop end point, packets drop in a smooth, linear pattern (called an interpolated graph) as shown in [Figure 11 on page 292](#). After the drop end point, all packets drop.

Figure 11: WRED-Drop Profile Packet Drop Pattern



The thick line in [Figure 11 on page 292](#) shows the packet drop characteristics for a sample WRED profile. At the drop start point, the queue reaches a fill level of 30 percent. At the drop end point, the queue fill level reaches 50 percent, and the maximum drop rate is 80 percent.

No packets drop until the queue fill level reaches the drop start point of 30 percent. When the queue reaches the 30 percent fill level, packets begin to drop. As the queue fills, the percentage of packets dropped increases in a linear fashion. When the queue fills to the drop end point of 50 percent, the rate of packet drop has increased to the maximum drop rate of 80 percent. When the queue fill level exceeds the drop end point of 50 percent, all of the packets drop until the queue fill level drops below 50 percent.

Defining Drop Profiles on QFX10000 Switches

Each queue fill level pairs with a drop probability. As the queue fills to different levels, every time it reaches a fill level configured in a drop profile, the queue applies the drop probability paired with that fill level to the traffic in the queue that exceeds the fill level. You can configure up to 32 pairs of fill levels and drop probabilities to create a customized packet drop probability curve with up to 32 points of differentiation.

Packets are not dropped until they reach the first configured queue fill level. When the queue reaches the first fill level, packets begin to drop at the configured drop probability rate paired with the first fill level. When the queue reaches the second fill level, packets begin to drop at the configured drop probability rate paired with the second fill level. This process continues for the number of fill level/drop probability pairs that you configure in the drop profile.

Drop profiles are interpolated, not segmented. An interpolated drop profile gradually increases the drop probability along a curve between each configured fill level. When the queue reaches the next fill level, the drop probability reaches the drop probability paired with that fill level. A segmented drop profile “jumps” from one fill level and drop probability setting to another in a stepped fashion. The drop probability of traffic does not change as the queue fills until the next fill level is reached.

An example of interpolation is a drop profile with three fill level/drop probability pairs:

- 25 percent queue fill level paired with a 30 percent drop probability
- 50 percent queue fill level paired with a 60 percent drop probability
- 75 percent queue fill level paired with a 100 percent drop probability (all packets that exceed the 75 percent queue fill level are dropped)

The queue drops no packets until its fill level reaches 25 percent. During periods of congestion, when the queue fills above 25 percent full, the queue begins to drop packets at a rate of 30 percent of the packets above the fill level.

However, as the queue continues to fill, it does not continue to drop packets at the 30 percent drop probability. Instead, the drop probability gradually increases as the queue fills to the 50 percent fullness level. When the queue reaches the 50 percent fill level, the drop probability has increased to the configured drop probability pair for the fill level, which is 60 percent.

As the queue continues to fill, the drop probability does not remain at 60 percent, but continues to rise as the queue fills. When the queue reaches the final fill level at 75 percent full, the drop probability has risen to 100 percent and all packets that exceed the 75 percent fill level are dropped.

Default Drop Profile

If you do not configure drop profiles and apply them to queue schedulers, the device uses the default drop profile for lossy traffic classes. In the default drop profile, when the fill level is 0 percent, the drop probability is 0 percent. When the fill level is 100 percent, the drop probability is 100 percent. During periods of congestion, as soon as packets arrive on a queue, the default profile might begin to drop packets.

Packet Drop Method

When a packet reaches the head of a queue, the device calculates a random number between 0 and 100. The device plots the random number against the drop profile using the current fill level of the queue. When the random number falls above the graph line, the queue transmits the packet out the egress interface. When the number falls below graph the line, the device drops the packet.

Packet Drop Example for Switches Except QFX10000

To create the linear drop pattern from the drop start point to the drop end point, the drop probabilities are derived using a linear approximation with eight sections, or steps, from the minimum queue fill level to the maximum queue fill level. The fill levels are divided into the eight sections equally, starting at the minimum fill level and ending at the maximum fill level. As the queue fills, the percentage of dropped packets increases. The percentage of packets dropped is based on the maximum drop rate.

For example, the default drop profile (which specifies a maximum drop rate of 100 percent) has the following drop probabilities at each section, or step, in the eight-section linear drop pattern:

- First section—The minimum drop probability is 6.25 percent of the maximum drop rate. The maximum drop probability is 12.5 percent of the maximum drop rate.
- Second section—The minimum drop probability is 18.75 percent of the maximum drop rate. The maximum drop probability is 25 percent of the maximum drop rate.
- Third section—The minimum drop probability is 30.25 percent of the maximum drop rate. The maximum drop probability is 37.5 percent of the maximum drop rate.
- Fourth section—The minimum drop probability is 43.75 percent of the maximum drop rate. The maximum drop probability is 50 percent of the maximum drop rate.
- Fifth section—The minimum drop probability is 56.25 percent of the maximum drop rate. The maximum drop probability is 62 percent of the maximum drop rate.
- Sixth section—The minimum drop probability is 68.75 percent of the maximum drop rate. The maximum drop probability is 75.5 percent of the maximum drop rate.
- Seventh section—The minimum drop probability is 81.25 percent of the maximum drop rate. The maximum drop probability is 87.5 percent of the maximum drop rate.
- Eighth section—The minimum drop probability is 92.75 percent of the maximum drop rate. The maximum drop probability is 100 percent of the maximum drop rate.

Packets drop even when there is no congestion, because packet drops begin at the drop start point regardless of whether congestion exists on the port. The default drop profile example represents the worst-case scenario, because the drop start point fill level is 0 percent, so packet drop begins when the queue starts to receive packets.

You can specify when packets begin to drop by configuring a drop start point at a fill level greater than 0 percent. For example, if you configure a drop profile that has a drop start point of 30 percent, packets do not drop until the queue is 30 percent full. We recommend that you configure drop profiles that are appropriate to your network traffic conditions.

The smaller the gap between the minimum drop rate (which is always 0) and the maximum drop rate, the smaller the gap between the minimum drop probability and the maximum drop probability at each

section (step) of the linear drop pattern. The default drop profile, which has the maximum gap between the minimum drop rate (0 percent) and the maximum drop rate (100 percent), has the highest gap between the minimum drop probability and the maximum drop probability at each step. Configuring a lower maximum drop rate for a drop profile reduces the gap between the minimum drop probability and the maximum drop probability.

Drop Profile Maps

Drop profile maps are part of scheduler configuration. A drop profile map maps drop profiles to packet loss priorities. Specifying the drop profile map in a scheduler associates the drop profile with the forwarding classes (queues) that you map to the scheduler in a scheduler map.

You configure loss priority for a queue in the classifier section of the CoS configuration hierarchy, and the loss priority is applied to the traffic assigned to the forwarding class at the ingress interface.

Each scheduler can have multiple drop profile maps.

Congestion Prevention

Configuring drop profiles on output queues enables you to control how congestion affects other queues on a port. If you do not configure drop profiles and map them to output queues, the device uses the default drop profile on queues that forward lossy traffic.

For example, if an ingress port forwards traffic to more than one egress port, and at least one of the egress ports experiences congestion, that can cause ingress port congestion. Ingress port congestion (ingress buffer exceeds its resource allocation) can cause frames to drop at the ingress port instead of at the egress port. Ingress port frame drop affects all of the egress ports to which the congested ingress port forwards traffic, not just the congested egress port.

NOTE: Do not configure drop profiles for the `fcoe` and `no-loss` forwarding classes. FCoE and other lossless traffic queues require lossless behavior (traffic queues that are configured with the `no-loss` packet drop attribute). Use priority-based flow control (PFC) to prevent frame drop on lossless priorities.

Configuring a WRED Drop Profile and Applying it to an Output Queue

To configure a WRED packet drop profile and apply it to an output queue:

1. Configure a drop profile:

- On switches except QFX10000 use the statement `set class-of-service drop-profiles profile-name interpolate fill-level drop-start-point fill-level drop-end-point drop-probability 0 drop-probability percentage`.
 - On QFX10000 switches use the statement `set class-of-service drop-profiles profile-name interpolate fill-level level1 level2 ... level32 drop-probability probability1 probability2 ... probability32`. You can specify as few as two fill level/drop probability pairs or as many as 32 pairs.
2. Map the drop profile to a queue scheduler using the statement `set class-of-service schedulers scheduler-name drop-profile-map loss-priority (low | medium-high | high) protocol any drop-profile profile-name`. The name of the drop-profile is the name of the WRED profile configured in Step 1.
 3. Map the scheduler, which Step 2 associates with the drop profile, to the output queue using the statement `set class-of-service scheduler-maps map-name forwarding-class forwarding-class-name scheduler scheduler-name`. The forwarding class identifies the output queue. Forwarding classes are mapped to output queues by default, and can be remapped to different queues by explicit user configuration. The scheduler name is the scheduler configured in Step 2.
 4. On switches except QFX10000, associate the scheduler map with a traffic control profile using the statement `set class-of-service traffic-control-profiles tcp-name scheduler-map map-name`. The scheduler map name is the name configured in Step 3.
 5. On switches except QFX10000, associate the traffic control profile with an interface using the statement `set class-of-service interfaces interface-name forwarding-class-set forwarding-class-set-name output-traffic-control-profile tcp-name`. The output traffic control profile name is the name of the traffic control profile configured in Step 4.

The interface uses the scheduler map in the traffic control profile to apply the drop profile (and other attributes) to the output queue (forwarding class) on that interface. Because you can use different traffic control profiles to map different schedulers to different interfaces, the same queue number on different interfaces can handle traffic in different ways.

6. On QFX10000 switches, associate the scheduler map with an interface using the statement `set class-of-service interfaces interface-name scheduler-map scheduler-map-name`.

The interface uses the scheduler map to apply the drop profile (and other attributes) to the output queue mapped to the forwarding class on that interface. Because you can use different scheduler maps on different interfaces, the same queue number on different interfaces can handle traffic in different ways.

Drop Profiles on Explicit Congestion Notification Enabled Queues

You must configure a WRED drop profile on queues that you enable for explicit congestion notification (ECN). On ECN-enabled queues, the drop profile sets the threshold for when the queue should mark a packet as experiencing congestion (see [Understanding CoS Explicit Congestion Notification](#)). When a

queue fills to the level at which the WRED drop profile has a packet drop probability greater than zero (0), the device might mark a packet as experiencing congestion. Whether or not a device marks a packet as experiencing congestion is the same probability as the drop probability of the queue at that fill level.

On ECN-enabled queues, the device does not use the drop profile to control dropping packets that are not ECN-capable packets (packets marked non-ECT, ECN code bits 00) during periods of congestion. Instead, the device uses the tail-drop algorithm to drop non-ECN-capable packets during periods of congestion. When a queue fills to its maximum level of fullness, tail-drop simply drops all subsequently arriving packets until there is space in the queue to buffer more packets. All non-ECN-capable packets are treated the same way.

To apply a WRED drop profile to non-ECT traffic, configure a multifield (MF) classifier to assign non-ECT traffic to a different output queue that is not ECN-enabled, and then apply the WRED drop profile to that queue.

RELATED DOCUMENTATION

Understanding Junos CoS Components

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Understanding CoS Explicit Congestion Notification

Example: Configuring WRED Drop Profiles

Example: Configuring Drop Profile Maps

Example: Configuring Unicast Classifiers

Configuring CoS WRED Drop Profiles

Configuring CoS Drop Profile Maps

Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)

Configuring CoS WRED Drop Profiles

IN THIS SECTION

- [Drop Profiles on Switches Except QFX10000 | 298](#)
- [Drop Profiles on QFX 10000 Switches | 299](#)

You can configure an interpolated weighted random early detection (WRED) profile to control traffic congestion by controlling packet drop characteristics for different packet loss priorities.

Drop profiles specify two values, which work as pairs:

- **Fill level**—The queue fullness value, which represents a percentage of the memory used to store packets in relation to the total amount of memory allocated to the queue.
- **Drop probability**—The percentage value that corresponds to the likelihood that an individual packet is dropped.

NOTE: Do not enable WRED on lossless traffic flows (forwarding classes configured with the no-loss packet drop attribute). Use priority-based flow control (PFC) to prevent packet loss on lossless forwarding classes.

Except on QFX10000, you cannot enable WRED on multidestination (multicast) queues on. You can enable WRED only on unicast queues.

OCX Series switches do not support lossless flows or PFC.

NOTE: On ECN-enabled queues, the drop profile sets the threshold for when the queue should mark a packet as experiencing congestion (see [Understanding CoS Explicit Congestion Notification](#)). On ECN-enabled queues, the switch does not use the drop profile to control dropping packets that are not ECN-capable packets during periods of congestion. Instead, the switch uses the tail-drop algorithm to drop non-ECN-capable packets during periods of congestion. When a queue fills to its maximum level of fullness, tail-drop simply drops all subsequently arriving packets until there is space in the queue to buffer more packets. All non-ECN-capable packets are treated the same way.

Drop Profiles on Switches Except QFX10000

Interpolated means that the switch creates a smooth drop curve from a drop start point to a drop end point, with a maximum drop rate that is reached at the drop end point.

The dropstart point is the average queue fill level when the WRED algorithm starts to drop packets. Before the drop start point, no packets are scheduled to drop. Specify the drop start point using the first of two fill-level statements.

The drop end point is the average queue fill level at which all subsequently arriving packets are dropped. When the queue fill levels falls below the drop end point, packets begin to be forwarded again. (At the drop end point, the packet drop probability becomes 100 percent.) Specify the drop end point using the second of two fill-level statements.

The minimum drop rate is always 0. Specify the minimum drop rate using the first of two drop-probability statements. The maximum drop rate is the drop probability when the average queue fill level reaches the drop end point. Specify the maximum drop rate using the second of two drop-probability statements.

The drop rate is zero until the queue fill level reaches the drop start point. As the queue continues to fill, packets drop in smooth linear curve until the queue reaches the drop end point, when packets drop at the maximum drop rate. If the queue fills beyond the drop end point, all packets that match the drop profile are dropped.

To configure a WRED profile using the CLI on switches except QFX10000:

1. Name the drop profile and set the drop start point, drop end point, minimum drop rate, and maximum drop rate for the drop profile:

```
[edit class-of-service]
user@switch# set drop-profile drop-profile-name interpolate fill-level percentage fill-level
percentage drop-probability 0 drop-probability percentage
```

Drop Profiles on QFX 10000 Switches

Each queue fill level pairs with a drop probability. As the queue fills to different levels, every time it reaches a fill level configured in a drop profile, the queue applies the drop probability paired with that fill level to the traffic in the queue that exceeds the fill level. You can configure up to 32 pairs of fill levels and drop probabilities to create a customized packet drop probability curve with up to 32 points of differentiation.

Packets are not dropped until they reach the first configured queue fill level. When the queue reaches the first fill level, packets begin to drop at the configured drop probability rate paired with the first fill level. When the queue reaches the second fill level, packets begin to drop at the configured drop probability rate paired with the second fill level. This process continues for the number of fill level/drop probability pairs that you configure in the drop profile.

Drop profiles are *interpolated*. An interpolated drop profile gradually increases the drop probability along a curve between each configured fill level. When the queue reaches the next fill level, the drop probability reaches the drop probability paired with that fill level.

To configure a WRED profile using the CLI on QFX10000 switches:

1. Name the drop profile and set the fill levels and their associated drop probabilities as percentages. For every fill level, there must be a paired drop probability (you must configure the same number of fill levels and drop probabilities).

```
[edit class-of-service]
user@switch# set drop-profile drop-profile-name interpolate fill-level level1 level2 ...
level32 drop-probability probability1 probability2 ... probability32
```

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring WRED Drop Profiles

Defining CoS Queue Schedulers

Defining CoS Queue Schedulers for Port Scheduling

Configuring CoS Drop Profile Maps

Understanding CoS WRED Drop Profiles

Example: Configuring WRED Drop Profiles

IN THIS SECTION

- [Requirements | 301](#)
- [Overview | 301](#)
- [Configuring WRED Drop Profiles on Switches Except QFX10000 | 302](#)
- [Configuring WRED Drop Profiles on QFX10000 Switches | 305](#)

You can configure interpolated weighted random early detection (WRED) profiles to control traffic congestion by controlling packet drop characteristics for different packet loss priorities.

NOTE: Do not enable WRED on lossless traffic flows. Use priority-based flow control (PFC) to prevent packet loss on lossless forwarding classes. (OCX Series switches do not support lossless flows or PFC.)

Except on QFX10000 switches, you cannot enable WRED on multidestination (multicast) queues. You can enable WRED only on unicast queues.

Requirements

This example uses the following hardware and software components:

- One switch
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series or Junos OS Release 15.1X53-D10 or later for the QFX10000.

Overview

You associate WRED drop profiles with loss priorities in a scheduler. When you map the scheduler to a forwarding class (queue), you apply the interpolated drop profile to traffic of the specified loss priority on that queue. Drop profiles specify two values, which work as pairs:

- Fill level—The queue fullness value, which represents a percentage of the memory used to store packets in relation to the total amount of memory allocated to the queue.
- Drop probability—The percentage value that corresponds to the likelihood that an individual packet is dropped.

NOTE: On ECN-enabled queues, the drop profile sets the threshold for when the queue should mark a packet as experiencing congestion (see [Understanding CoS Explicit Congestion Notification](#)). On ECN-enabled queues, the switch does not use the drop profile to control dropping packets that are not ECN-capable packets during periods of congestion. Instead, the switch uses the tail-drop algorithm to drop non-ECN-capable packets during periods of congestion. When a queue fills to its maximum level of fullness, tail-drop simply drops all subsequently arriving packets until there is space in the queue to buffer more packets. All non-ECN-capable packets are treated the same way.

Configuring WRED Drop Profiles on Switches Except QFX10000

IN THIS SECTION

- [Verification](#) | 304

Configuration

Step-by-Step Procedure

Interpolated means that the switch creates a smooth drop curve from a drop start point to a drop end point, with a maximum drop rate that is reached at the drop end point:

- Drop start point—Percentage of average queue fill level when the WRED algorithm starts to drop packets. Before the drop start point, no packets are scheduled to drop.
- Drop end point—Average queue fill level at which all subsequently arriving packets are dropped. When the queue fill levels falls below the drop end point, packets begin to be forwarded again. (At the drop end point, the packet drop probability becomes 100 percent.)
- Maximum drop rate—Drop probability when the average queue fill level reaches the drop end point.

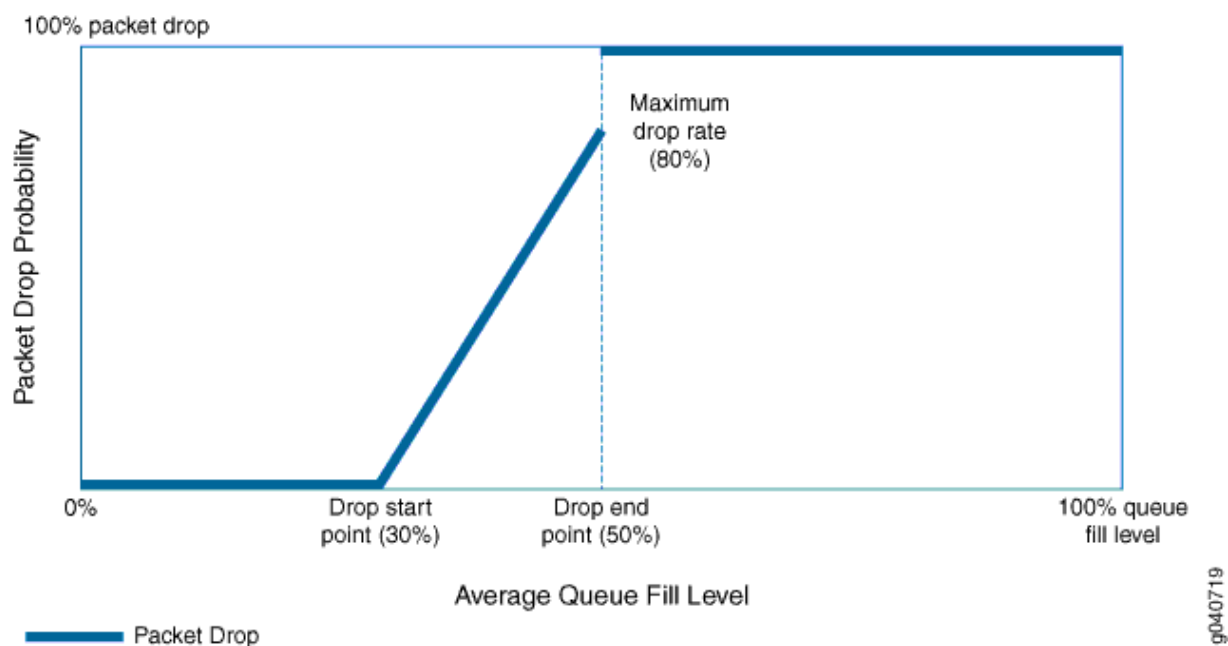
You set the drop start point and the drop end point by specifying two queue fill level percentage values. The first value is the drop start point and the second value is the drop end point.

You set the maximum drop rate by specifying two drop probability percentage values. The first value is always zero (0), which is the minimum drop rate, the probability of dropping a packet at the drop start point. The second value is the maximum drop rate at the drop end point.

The drop rate is zero until the queue fill level reaches the drop start point. As the queue continues to fill, packets drop in smooth linear curve until the queue reaches the drop end point, when packets drop at the maximum drop rate. If the queue fills beyond the drop end point, all packets that match the drop profile are dropped.

[Figure 12 on page 303](#) shows the graph for a drop profile with a drop start point of 30 percent, a drop end point of 50 percent, and a maximum drop rate of 80 percent.

Figure 12: WRED Drop Profile Packet Drop Example



The graph shows that when the queue fill level is less than 30 percent, the packet drop rate is zero. When the queue fill level reaches 30 percent, packets begin to drop. As the queue fills, a higher percentage of packets drop. When the queue fill level reaches 50 percent, the packet drop rate has climbed to 80 percent. When the queue fill level exceeds 50 percent, all packets drop.

This example describes how to configure the drop profile shown in [Figure 12 on page 303](#). The drop profile will have:

- The name be-dp1
- 30 percent for the drop start point (first fill-level setting)
- 50 percent for the drop end point (second fill-level setting)
- 0 percent for the minimum drop rate (first drop-probability setting)
- 80 percent for the maximum drop rate (second drop-probability setting)

You apply a drop profile by configuring a drop profile map that maps the drop profile to a packet loss priority, and associate the drop profile and packet loss priority with a scheduler. When you map the scheduler to a forwarding class (queue), the switch applies the drop profile to the packets in the forwarding class that have a matching packet loss priority.

1. Set the drop start point at 30 percent, the drop end point at 50 percent, the minimum drop rate at 0 percent, and the maximum drop rate at 80 percent for the drop profile be-dp1:

```
[edit class-of-service]
user@switch# set drop-profile be-dp1 interpolate fill-level 30 fill-level 50 drop-probability
0 drop-probability 80
```

Verification

IN THIS SECTION

- [Verifying the Drop Profile Configuration | 304](#)

Verifying the Drop Profile Configuration

Purpose

Verify that you configured the drop profile be-dp1 with the correct drop start and end points and with the correct drop rates.

Action

Verify the results of the drop profile configuration using the operational mode command `show configuration class-of-service drop-profiles be-dp1`:

```
user@switch> show configuration class-of-service drop-profiles be-dp1
interpolate {
    fill-level [ 30 50 ];
    drop-probability [ 0 80 ];
}
```

Configuring WRED Drop Profiles on QFX10000 Switches

IN THIS SECTION

- [Verification | 306](#)

Configuration

Step-by-Step Procedure

Each queue fill level pairs with a drop probability. As the queue fills to different levels, every time it reaches a fill level configured in a drop profile, the queue applies the drop probability paired with that fill level to the traffic in the queue that exceeds the fill level. You can configure up to 32 pairs of fill levels and drop probabilities to create a customized packet drop probability curve with up to 32 points of differentiation.

Packets are not dropped until they reach the first configured queue fill level. When the queue reaches the first fill level, packets begin to drop at the configured drop probability rate paired with the first fill level. When the queue reaches the second fill level, packets begin to drop at the configured drop probability rate paired with the second fill level. This process continues for the number of fill level/drop probability pairs that you configure in the drop profile.

Drop profiles are *interpolated*. An interpolated drop profile gradually increases the drop probability along a curve between each configured fill level. When the queue reaches the next fill level, the drop probability reaches the drop probability paired with that fill level.

This example describes how to configure a drop profile with three fill level/drop probability pairs:

- Drop profile name—be-dp1
- Queue fill levels—25 percent, 50 percent, 75 percent
- Drop probabilities—30 percent, 60 percent, 100 percent

Each of the three fill levels pairs with a drop probability to program the interpolated drop profile curve.

You apply a drop profile by configuring a drop profile map that maps the drop profile to a packet loss priority, and associate the drop profile and packet loss priority with a scheduler. When you map the scheduler to a forwarding class (queue), the switch applies the drop profile to the packets in the forwarding class that have a matching packet loss priority.

To configure a drop profile:

1. Set the drop start point at a 25 percent fill level, an intermediate fill level of 50 percent, and a drop end point of 75 percent. Set the paired drop probabilities to 30 percent, 60 percent, and 100 percent, respectively, for drop profile be-dp1:

```
[edit class-of-service]
user@switch# set drop-profile be-dp1 interpolate fill-level [ 25 50 75 ] drop-probability
[ 30 60 100 ]
```

Verification

IN THIS SECTION

- [Verifying the Drop Profile Configuration | 306](#)

Verifying the Drop Profile Configuration

Purpose

Verify that you configured the drop profile be-dp1 with the correct fill levels and drop probabilities.

Action

Verify the results of the drop profile configuration using the operational mode command `show configuration class-of-service drop-profiles be-dp1`:

```
user@switch> show configuration class-of-service drop-profiles be-dp1
interpolate {
    fill-level [ 25 50 75 ];
    drop-probability [ 30 60 100 ];
}
```

Configuring CoS Drop Profile Maps

A drop-profile map associates weighted random early detection (WRED) profiles for traffic of specified packet loss priorities with a scheduler. When you use a scheduler map to map a scheduler to a forwarding class, the drop profile map associated with the scheduler applies the specified WRED drop profile to traffic in the forwarding class that matches the specified packet loss priority.

Drop profile maps enable you to configure different drop profiles for traffic of different packet loss priorities within the same scheduler. You can associate different drop profiles with low-priority, medium-high priority, and high-priority traffic within a single scheduler, and then map that scheduler to a forwarding class. This applies the appropriate drop profile to traffic of each loss priority in a forwarding class. Drop profile maps apply to all traffic protocols.

To configure a drop-profile map:

- For the desired scheduler, configure the traffic loss priority and specify the drop profile you want to use to control the drop characteristics for traffic of that loss priority:

```
[edit class-of-service]
user@switch# set schedulers scheduler-name drop-profile-map loss-priority level protocol any
drop-profile drop-profile-name
```

NOTE: QFX10000 switches do not support the `protocol any` portion of the configuration. Drop profiles apply to all protocols.

Example: Configuring Drop Profile Maps

IN THIS SECTION

- [Requirements | 309](#)
- [Overview | 309](#)
- [Verification | 309](#)

A drop-profile map associates weighted random early detection (WRED) profiles for traffic of specified packet loss priorities with a scheduler. When you use a scheduler map to map a scheduler to a forwarding class, the drop profile map associated with the scheduler applies the specified WRED drop profile to traffic in the forwarding class that matches the specified packet loss priority.

Configuring a Drop Profile Map

CLI Quick Configuration

To quickly configure a drop profile map, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level.

```
[edit class-of-service]
set schedulers mylan drop-profile-map loss-priority low protocol any drop-profile lp-profile
set schedulers mylan drop-profile-map loss-priority medium-high protocol any drop-profile mh-profile
set schedulers mylan drop-profile-map loss-priority high protocol any drop-profile h-profile
```

Step-by-Step Procedure

To configure a drop profile map:

1. Configure the drop profile for low-priority traffic:

```
[edit class-of-service]
user@switch# set schedulers mylan drop-profile-map loss-priority low protocol any drop-profile lp-profile
```

2. Configure the drop profile for medium-high priority traffic:

```
[edit class-of-service]
user@switch# set schedulers mylan drop-profile-map loss-priority medium-high protocol any drop-profile mh-profile
```

3. Configure the drop profile for high-priority traffic:

```
[edit class-of-service]
user@switch# set schedulers mylan drop-profile-map loss-priority high protocol any drop-
profile h-profile
```

Requirements

This example uses the following hardware and software components:

- A Juniper Networks QFX3500 Switch
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series.

Overview

Drop profile maps enable you to configure different drop profiles for traffic of different packet loss priorities within the same scheduler. You can associate different drop profiles with low-priority, medium-high priority, and high-priority traffic within a single scheduler, and then map that scheduler to a forwarding class. This applies the appropriate drop profile to traffic of each loss priority in a forwarding class. Drop profile maps apply to all traffic protocols.

The following example describes how to configure a drop profile map for a scheduler named `mylan` that includes:

- A drop profile called `lp-profile` for low-priority traffic
- A drop profile called `mh-profile` for medium-high priority traffic
- A drop profile called `h-profile` for high-priority traffic

You apply the drop profiles in the drop profile map to a forwarding class by associating the scheduler `mylan` with a forwarding class in a scheduler map.

Verification

IN THIS SECTION

- [Verifying the Drop Profile Map Configuration | 310](#)

Verifying the Drop Profile Map Configuration

Purpose

Verify that you configured the drop profile map for the scheduler `mylan` with the correct loss priorities and drop profiles.

Action

Verify the results of the drop profile map configuration using the operational mode command `show configuration class-of-service schedulers mylan`:

```
user@switch> show configuration class-of-service schedulers mylan
transmit-rate 3g;
shaping-rate percent 100;
priority low;
drop-profile-map loss-priority low protocol any drop-profile lp-profile;
drop-profile-map loss-priority medium-high protocol any drop-profile mh-profile;
drop-profile-map loss-priority high protocol any drop-profile h-profile;
```

NOTE: This example does not include configuring scheduler bandwidth and priority. This information (transmit rate, shaping rate, and priority) is shown for completeness.

RELATED DOCUMENTATION

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

Example: Configuring Queue Schedulers

Example: Configuring Queue Schedulers for Port Scheduling

Example: Configuring WRED Drop Profiles

Configuring CoS Drop Profile Maps

Understanding CoS WRED Drop Profiles

Troubleshooting a Port Reset on QFabric Systems When a Queue Stops Transmitting Traffic

IN THIS SECTION

- [Problem | 311](#)
- [Cause | 311](#)
- [Solution | 311](#)

Problem

Description

In QFabric systems, if any queue that contains outgoing packets does not transmit packets for 12 consecutive seconds, the port automatically resets.

Cause

Failure of a queue to transmit packets for 12 consecutive seconds may be due to:

- A strict-high priority queue consuming all of the port bandwidth
- Several queues consuming all of the port bandwidth
- Any queue or port receiving continuous priority-based flow control (PFC) or 802.3x Ethernet PAUSE messages (received PFC and PAUSE messages prevent a queue or a port, respectively, from transmitting packets because of network congestion)
- Other conditions that prevent a queue from obtaining port bandwidth for 12 consecutive seconds

Solution

If the cause is a strict-high priority queue or other queues consuming all of the port bandwidth, you can use rate shaping to configure a maximum rate for the queues that are using all of the port bandwidth and preventing other queues from obtaining bandwidth on the port. You configure a maximum rate by creating a scheduler, using a scheduler map to apply it to a forwarding class (which maps to an output queue), and applying the scheduler map to the port using a forwarding class set and a traffic control profile.

To configure rate shaping using the CLI:

1. Name the existing scheduler or create a scheduler and define the maximum bandwidth as a rate or as a percentage:

```
[edit class-of-service]
user@switch# set schedulers scheduler-name shaping-rate (rate | percent
percentage)
```

2. Configure a scheduler map to associate the scheduler with the forwarding class (queue) that is consuming all of the port bandwidth:

```
[edit class-of-service]
user@switch# set scheduler-maps scheduler-map-name forwarding-class forwarding-class-name
scheduler scheduler-name
```

3. Associate the scheduler map with a traffic control profile:

```
[edit class-of-service]
user@switch# set traffic-control-profiles traffic-control-profile-name scheduler-map
scheduler-map-name
```

4. Associate the traffic control profile (and thus the scheduler map that contains the rate shaping queue scheduler) with a forwarding class set and apply them to the interface that is being reset:

```
[edit class-of-service]
user@switch# set interfaces interface-name forwarding-class-set fc-set-name output-traffic-
control-profile traffic-control-profile-name
```

For example, a strict-high priority queue is using all of the bandwidth on interface `shpnode:xe-0/0/10` and preventing other queues from transmitting for 12 consecutive seconds. You decide to set a maximum rate of 7 Gbps on the strict-high priority queue to ensure that at least 3 Gbps of the port bandwidth is available to service other queues. [Table 65 on page 313](#) shows the topology for this example:

Table 65: Components of the Rate Shaping Troubleshooting Example

Component	Settings
Affected interface	shpnode:xe-0/0/10
Scheduler (strict-high priority scheduler)	Name: shp-sched Shaping rate: 7g Priority: strict-high NOTE: This example assumes that the scheduler already exists and has been configured as strict-high priority, but that rate shaping to prevent the strict-high priority traffic from using all of the port bandwidth has not been applied.
Scheduler map	Name: shp-map Forwarding class to associate with the shp-sched scheduler: strict-high NOTE: This example assumes that a strict-high priority forwarding class has been configured and assigned the name strict-high.
Traffic control profile	Name: shp-tcp NOTE: This example does not describe how to define a complete traffic control profile.
Forwarding class set	Name: shp-pg

To configure the scheduler, map it to the strict-high priority forwarding class, and apply it to interface shpnode:xe-0/0/10 using the CLI:

1. Specify the scheduler for the strict-high priority queue (shp-sched) with a maximum bandwidth of 7 Gbps:

```
[edit class-of-service schedulers]
user@switch# set shp-sched shaping-rate 7g
```

2. Configure a scheduler map (shp-map) that associates the scheduler (shp-sched) with the forwarding class (strict-high):

```
[edit class-of-service scheduler-maps]
user@switch# set shp-map forwarding-class strict-high scheduler shp-sched
```

3. Associate the scheduler map shp-map with a traffic control profile (shp-tcp):

```
[edit class-of-service traffic-control-profiles]
user@switch# set shp-tcp scheduler-map shp-map
```

4. Associate the traffic control profile shp-tcp with a forwarding class set (shp-pg) and the affected interface (shpnode:xe-0/0/10):

```
[edit class-of-service]
user@switch# set interfaces shpnode:xe-0/0/10 forwarding-class-set shp-pg output-traffic-
control-profile shp-tcp
```

RELATED DOCUMENTATION

[Understanding CoS Output Queue Schedulers](#)

[Defining CoS Queue Scheduling Priority](#)

[Example: Configuring Queue Schedulers](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\)](#)

[Example: Configuring Forwarding Class Sets](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#)

Using Schedulers (Interconnect Device Fabric)

IN THIS CHAPTER

- Understanding Default CoS Scheduling on QFabric System Interconnect Devices (Junos OS Release 13.1 and Later Releases) | **315**
- Understanding CoS Scheduling Across the QFabric System | **327**
- Example: Configuring CoS Scheduling Across the QFabric System | **353**
- Understanding CoS Fabric Forwarding Class Sets | **396**
- Configuring CoS Fabric Forwarding Class Set Scheduler Maps (Fabric Scheduler to Fabric FC-Set Mapping) | **409**
- Understanding How to Mitigate Fate Sharing on a QFabric System Interconnect Device by Remapping Traffic Flows (Forwarding Classes) | **410**
- Configuring Fate Sharing Mitigation Across the Interconnect Device by Remapping Traffic Flows (Forwarding Classes) | **434**

Understanding Default CoS Scheduling on QFabric System Interconnect Devices (Junos OS Release 13.1 and Later Releases)

IN THIS SECTION

- Hierarchical CoS Architecture Across a QFabric System Interconnect Device | **316**
- Default CoS on Interconnect Device Fabric Interfaces | **317**

The default class-of-service (CoS) properties on the QFabric system Interconnect device interfaces are optimized to best utilize the fabric resources. You cannot configure CoS properties on QFabric System Interconnect device interfaces.

Hierarchical CoS Architecture Across a QFabric System Interconnect Device

Because Interconnect devices support traffic from multiple Node devices that have multiple CoS configurations, CoS on Interconnect device fabric interfaces differs from CoS on Node device access and fabric interfaces.

The hierarchical CoS scheduling structure on the Interconnect device interfaces consists of two tiers:

1. Fabric forwarding class sets—Similar to fc-sets on Node devices, fabric fc-sets group traffic for transport across the Interconnect device fabric. Fabric fc-sets are global and apply to all traffic that crosses the fabric from all Node devices. See [Understanding CoS Fabric Forwarding Class Sets](#) for a detailed description of fabric fc-sets.
2. Class groups—Fabric fc-sets are grouped into class groups for transport across the Interconnect device.

Node devices and Interconnect devices each have a two-tier hierarchical CoS scheduling architecture. The architectures are slightly different, but each tier of the scheduling hierarchy performs analogous functions, as shown in [Table 66 on page 316](#).

Table 66: Hierarchical Scheduler Architecture on Node Devices and Interconnect Devices

Bandwidth Pool	Bandwidth Configuration on Node Devices	Bandwidth Configuration on Interconnect Devices
Port—Entire amount of bandwidth available to traffic on a port.	Access (xe) or fabric (fte) interfaces	Fabric (fte) or Clos fabric (bfte) interfaces
Priority group—Group of traffic types that requires similar CoS treatment. Each priority group receives a portion of the total available port bandwidth.	Forwarding class set (fc-set)	Class group
Priority—Most granular tier of bandwidth allocation. Each priority receives a portion of the total available priority group bandwidth.	Forwarding class (mapped to output queue)	Fabric fc-set (mapped to output queue)

Fabric FC-Sets

Fabric fc-sets are groups of forwarding classes that receive similar CoS treatment across the Interconnect device. Fabric fc-sets are global to the QFabric system and apply to all traffic that traverses

the fabric, from all connected Node devices. The CoS on a fabric fc-set applies to all the traffic that belongs to that fabric fc-set.

For example, a fabric fc-set that includes the best-effort forwarding class handles all of the best-effort traffic from all of the connected Node devices that traverses the Interconnect device fabric.

There are 12 default fabric fc-sets, including 5 visible fabric fc-sets and 7 hidden fabric fc-sets. The five visible fabric fc-sets have forwarding classes mapped to them by default. By default, the seven hidden fabric fc-sets do not carry traffic, but you can map forwarding classes to the hidden fabric fc-sets if you want to use them.

You can configure the forwarding class membership of each fabric fc-set. However, you cannot create new fabric fc-sets, and you cannot delete the 12 default fabric fc-sets.

Each fabric fc-set is mapped to an output queue. Each fabric interface has 12 output queues, one for each of the 12 fabric fc-sets. The traffic from all of the forwarding classes mapped to a fabric fc-set uses that fabric fc-set's output queue.

Fabric fc-sets are grouped into class groups for transport across the Interconnect device.

Class Groups for Fabric FC-Sets

To transport traffic across the fabric, the fabric organizes the fabric fc-sets into three classes called *class groups*. Class groups are not user-configurable. The three class groups are:

- **Strict-high priority**—All traffic in the fabric fc-set `fabric_fcset_strict_high`. This class group includes the traffic in strict-high priority and network-control forwarding classes, and in any forwarding classes you create on a Node device that consist of strict-high priority traffic.
- **Unicast**—All traffic in the fabric fc-sets `fabric_fcset_be`, `fabric_fcset_noloss1`, and `fabric_fcset_noloss2`. This class group includes the traffic in the best-effort, `fcoe`, and no-loss forwarding classes, and the traffic in any forwarding classes you create on a Node device that consist of best-effort or lossless unicast traffic. If you use any of the hidden no loss fabric fc-sets (`fabric_fcset_noloss3`, `fabric_fcset_noloss4`, `fabric_fcset_noloss5`, or `fabric_fcset_noloss6`), that traffic is part of this class group.
- **Multidestination**—All traffic in the fabric fc-set `fabric_fcset_multicast1`. This class group includes the traffic in the `mcast` forwarding class and in any forwarding classes you create on a Node device that consist of multidestination traffic. If you use any of the hidden multidestination fabric fc-sets (`fabric_fcset_multicast2`, `fabric_fcset_multicast3`, or `fabric_fcset_multicast4`), that traffic is part of this class group.

Default CoS on Interconnect Device Fabric Interfaces

The Interconnect device interfaces use the default CoS configuration as described in these sections:

Default Class Group Scheduling

Default class group bandwidth scheduling is analogous to default fc-set (priority group) scheduling on a Node device. Default class group scheduling uses weighted round-robin (WRR) scheduling, in which each class group receives a portion of the total available fabric interface bandwidth, based on the class group's traffic type, as shown in [Table 67 on page 318](#):

Table 67: Class Group Default Scheduling Properties and Membership

Class Group	Fabric fc-sets	Forwarding Classes (Default Mapping)	Class Group Scheduling Properties (Weight)
Strict-high priority	fabric_fcset_strict_high	<ul style="list-style-type: none"> All strict-high priority forwarding classes network-control 	Traffic in the strict-high priority class group is served first. This class group receives all of the bandwidth it needs to empty its queues and therefore can starve other types of traffic during periods of high-volume strict priority traffic. Plan carefully and use caution when determining how much traffic to configure as strict-high priority traffic.

Table 67: Class Group Default Scheduling Properties and Membership *(Continued)*

Class Group	Fabric fc-sets	Forwarding Classes (Default Mapping)	Class Group Scheduling Properties (Weight)
Unicast	<ul style="list-style-type: none"> • fabric_fcset_be • fabric_fcset_nolos s1 • fabric_fcset_nolos s2 <p>Includes the hidden lossless fabric fc-sets if used:</p> <ul style="list-style-type: none"> • fabric_fcset_nolos s3 • fabric_fcset_nolos s4 • fabric_fcset_nolos s5 • fabric_fcset_nolos s6 	<ul style="list-style-type: none"> • best-effort • fcoe • no-loss <p>NOTE: No forwarding classes are mapped to the hidden lossless fabric_fcsets by default.</p>	Traffic in the unicast class group receives an 80% weight in the weighted round-robin (WRR) calculations. After the strict-high priority class group has been served, the unicast class group receives 80% of the remaining fabric bandwidth. (If more bandwidth is available, the unicast class group can use more bandwidth.)
Multidestination	<p>fabric_fcset_multicast1</p> <p>Includes the hidden multidestination fabric fc-sets if used:</p> <ul style="list-style-type: none"> • fabric_fcset_multicast2 • fabric_fcset_multicast3 • fabric_fcset_multicast4 	<ul style="list-style-type: none"> • mcast <p>NOTE: No forwarding classes are mapped to the hidden multidestination fabric_fcsets by default.</p>	Traffic in the multidestination class group receives a 20% weight in the WRR calculations. After the strict-high priority class group has been served, the multidestination class group receives 20% of the remaining fabric bandwidth. (If more bandwidth is available, the multidestination class group can use more bandwidth.)

Only the five visible fabric fc-sets have traffic mapped to them by default. The fabric fc-sets within each class group are weighted by their transmit rates (guaranteed minimum bandwidth), and they receive bandwidth from the class group's total bandwidth using weighted round-robin (WRR) scheduling.

Default Fabric FC-Set Scheduling

Default fabric fc-set bandwidth scheduling is analogous to default forwarding class (priority) scheduling on a Node device. Each fabric fc-set receives a guaranteed minimum percentage of the port bandwidth that the class group receives. The guaranteed minimum percentage is called the *transmit rate*.

[Table 68 on page 320](#) shows the default transmit rate for each of the default fabric fc-sets.

Table 68: Default Fabric FC-Set Scheduler Configuration

Default Fabric FC-Set	Transmit Rate (Percentage of Class Group Bandwidth)
fabric_fcset_strict_high	N/A Strict-high priority traffic is served first, before any other traffic is served. Strict-high priority traffic receives all of the bandwidth it needs to empty its queues and therefore can starve other types of traffic during periods of high-volume strict priority traffic. Plan carefully and use caution when determining how much traffic to configure as strict-high priority traffic.
fabric_fcset_noloss1	35%
fabric_fcset_noloss2	35%
fabric_fcset_be	10%
fabric_fcset_multicast1	20%

Each fabric fc-set belongs to a class group. Each class group receives a portion of the total available port bandwidth. Each fabric fc-set in a class group receives a portion of the total available class group bandwidth based on the transmit rate (weight) of the fabric fc-set.

Traffic in fabric_fcset_strict_high does not have a default transmit rate because fabric_fcset_strict_high receives all of the bandwidth needed to empty its queue before other queues are served. Traffic in the remaining fabric fc-sets receives bandwidth in a ratio proportional to the default transmit rate of each fabric fc-set.

Each of the following hidden fabric fc-sets receives a default scheduling weight of 1:

- fabric_fcset_noloss3
- fabric_fcset_noloss4
- fabric_fcset_noloss5
- fabric_fcset_noloss6
- fabric_fcset_multicast2
- fabric_fcset_multicast3
- fabric_fcset_multicast4

You must explicitly map forwarding classes to hidden fabric fc-sets if you want to use the hidden fabric fc-sets.



CAUTION: Bandwidth is allocated to fabric fc-sets based on scheduling weight. The scheduling weights of the visible (default) fabric fc-sets are the same as their transmit rates, so in the unicast class group, fabric_fcset_noloss1 and fabric_fcset_noloss2 each have a weight of 35 and fabric_fcset_be has a weight of 10. In the multidestination class group, the default fabric_fcset_multicast1 has a weight of 20. The hidden multicast and noloss fabric fc-sets each have a scheduling weight of 1.

The scheduling weights mean that when the visible fabric fc-sets are fully utilizing their allocated bandwidth:

- The hidden noloss fc-sets (fabric_fcset_noloss3, fabric_fcset_noloss4, fabric_fcset_noloss5, and fabric_fcset_noloss6) receive bandwidth at a proportional rate of 1:35 compared to the default noloss fc-sets.
- The hidden multicast fc-sets (fabric_fcset_multicast2, fabric_fcset_multicast3, and fabric_fcset_multicast4) receive bandwidth at a proportional rate of 1:20 compared to the default multicast fc-sets.

If you map traffic to a hidden fabric fc-set, that fabric fc-set receives the proportional amount of class group bandwidth that corresponds to its scheduling weight (1). The amount of bandwidth allocated to a hidden fabric fc-set depends on how much bandwidth the other fc-sets in the same class group consume. When the visible fabric fc-sets fully utilize their bandwidth, hidden fabric fc-sets receive only their minimum weight in bandwidth. (However, even a low scheduling weight results in a relatively large absolute bandwidth allocation because each fabric port is a 40-Gbps port.)

For example, if fabric_fcset_noloss1 and fabric_fcset_noloss2 each consume all of the 35 percent of bandwidth allocated to them, and fabric_fcset_be consumes all of the

10 percent of bandwidth allocated to it, then `fabric_fcset_noloss3`, `fabric_fcset_noloss4`, `fabric_fcset_noloss5`, and `fabric_fcset_noloss6` receive bandwidth at a rate of 1:80 compared to the visible noloss fabric fc-sets. (If the visible fabric fc-sets do not use all of their allocated bandwidth, then the hidden fabric fc-sets receive more bandwidth.)

Another example is if we map lossless traffic to `fabric_fcset_noloss3` and to `fabric_fcset_noloss4`. `Fabric_fcset_noloss1` uses 10 percent of its 35 percent allocation of unicast class group bandwidth. `Fabric_fcset_noloss2` uses 15 percent of its 35 percent allocation of unicast class group bandwidth. `Fabric_fcset_be` uses 5 percent of its allocated bandwidth. `Fabric_fcset_noloss3` and `fabric_fcset_noloss4` can use the remaining unicast class group bandwidth allocated to lossless traffic. However, if the traffic on `fabric_fcset_noloss1`, `fabric_fcset_noloss2`, or `fabric_fcset_be` increases, the bandwidth allocated to the hidden fabric fc-sets decreases.

Similarly, if you map traffic to a hidden multidestination fabric fc-set (`fabric_fcset_multicast2`, `fabric_fcset_multicast3`, `fabric_fcset_multicast4`), that multidestination fabric fc-set receives the proportional amount of class group bandwidth that corresponds to its scheduling weight (1). The amount of bandwidth allocated to a hidden multidestination fabric fc-set depends on how much bandwidth the other fc-sets in the multidestination class group consume. When `fabric_fcset_multicast1` (the visible fabric fc-set) fully utilizes its bandwidth, hidden fabric fc-sets receive only their minimum weight in bandwidth. For example, if `fabric_fcset_multicast1` uses its full bandwidth allocation, then the hidden multidestination fabric fc-sets receive bandwidth at a rate of 1:20 compared to `fabric_fcset_multicast1`.

Default Class Group and Fabric FC-Set Scheduling Example

The following example shows how default scheduling allocates the total port bandwidth among the class groups and their fabric fc-sets. In the example, traffic is mapped to each of the forwarding classes in the five visible fabric fc-sets, and the strict-high priority class group consumes an average of 10 percent of the 40-Gbps fabric interface bandwidth (4 gigabits), leaving 90 percent of the fabric interface bandwidth (36 gigabits) for the remaining class groups.

In this scenario, by default, the strict-high priority class group includes one fabric fc-set (`fabric_fcset_strict_high`), the unicast class group includes three fabric fc-sets (`fabric_fcset_be`, `fabric_fcset_noloss1`, and `fabric_fcset_noloss2`), and the multidestination class group includes one fabric fc-set (`fabric_fcset_multicast1`). Each individual fabric fc-set receives the following treatment:

- Strict-high priority class group (`fabric_fcset_strict_high`)—This group is assumed to average 10 percent (4 gigabits) for the purposes of this example. Because the strict-high priority class group is served first and receives all of the bandwidth it requires to empty its queue, in real networks the amount of

required bandwidth fluctuates and affects the amount of bandwidth available to the other class groups.

TIP: To prevent strict-high priority traffic from using too much bandwidth, you can set a maximum bandwidth limit by configuring a scheduler shaping rate for the `fabric_fcset_strict_high` fabric fc-set.

- Unicast class group (`fabric_fcset_be`, `fabric_fcset_noloss1`, and `fabric_fcset_noloss2`)—Each of these fabric fc-sets receives a weighted portion of the 80 percent of the total port bandwidth available after the strict-high traffic has been served. The weight corresponds to the transmit rate of each fabric fc-set. The following calculations show the minimum port bandwidth allocated to each of the unicast class group fabric fc-sets:

- `fabric_fcset_be`

$10 / (35 + 35 + 10)\%$ of 80% of the available port bandwidth (12.5 percent of 80 percent of port bandwidth)

The 10 that is the numerator in $10 / (35 + 35 + 10)$ is the percentage of bandwidth allocated to the `fabric_fcset_be` by the transmit rate weight. The $(35 + 35 + 10)$ in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the three fabric fc-sets in the unicast class group.

The 80 percent represents 80 percent of the port bandwidth available after strict-high priority traffic is served (36 gigabits).

The resulting equation is:

$10 / (35 + 35 + 10)\% \times (0.8 \times 36 \text{ gigabits}) = \text{approximately } 3.6 \text{ gigabits}$

- `fabric_fcset_noloss1` and `fabric_fcset_noloss2`

The default minimum bandwidth for the two visible lossless fabric fc-sets is the same because both of these fabric fc-sets have the same transmit rate weight.

$35 / (35 + 35 + 10)\%$ of 80% of the port bandwidth (43.75 percent of 80 percent of port bandwidth)

The 35 that is the numerator in $35 / (35 + 35 + 10)$ is the percentage of bandwidth allocated to each of the noloss fabric fc-sets by the transmit rate weight. The $(35 + 35 + 10)$ in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the three fabric fc-sets in the unicast class group.

The 80 percent represents 80 percent of the port bandwidth available after strict-high priority traffic is served (36 gigabits).

The resulting equation is:

$$35 / (35 + 35 + 10)\% \times (0.8 \times 36 \text{ gigabits}) = \text{approximately } 12.6 \text{ gigabits}$$

- **Multidestination class group (fabric_fcset_multicast1)**—Because only one fabric fc-set is configured by default in the multidestination class group, it receives 100 percent of the 20 percent of the total port bandwidth available to the multidestination class group after the strict-high traffic has been served:

$$100 / (100)\% \text{ of } 20\% \text{ of the available port bandwidth (100 percent of 20 percent of available port bandwidth)}$$

The resulting equation is:

$$100 / 100\% \times (0.2 \times 36 \text{ gigabits}) = \text{approximately } 7.2 \text{ gigabits}$$

Default PFC and Lossless Transport Across the Interconnect Device

The Interconnect device incorporates flow control mechanisms to support lossless transport during periods of congestion on the fabric. To support the priority-based flow control (PFC) feature on the Node devices, the Interconnect device fabric supports lossless transport for up to six IEEE 802.1p priorities when the following two configuration constraints are met:

1. The IEEE 802.1p priority used for the traffic that requires lossless transport is mapped to a lossless forwarding class (a forwarding class configured with the `no-loss` parameter or the default `fcoe` or `no-loss` forwarding class).
2. The lossless forwarding class must be mapped to one of the lossless fabric fc-sets (`fabric_fcset_noloss1`, `fabric_fcset_noloss2`, `fabric_fcset_noloss3`, `fabric_fcset_noloss4`, `fabric_fcset_noloss5`, or `fabric_fcset_noloss6`). If you do not explicitly map lossless forwarding classes to fabric fc-sets, lossless forwarding classes are mapped by default to lossless fabric fc-sets `fabric_fcset_noloss1` and `fabric_fcset_noloss2`.

When traffic meets these two constraints, the fabric propagates back-pressure from egress queues during periods of congestion. However, to achieve end-to-end lossless transport across the QFabric system, you must also configure a congestion notification profile to enable PFC on the Node device ingress interfaces. To achieve end-to-end lossless transport across the network, you must configure PFC on all of the devices in the lossless traffic path.

For all other combinations of IEEE 802.1p priority to forwarding class mapping and all other combinations of forwarding class to fabric fc-set mapping, the default congestion control mechanism is normal packet drop. For example:

- **Case 1**—If the IEEE 802.1p priority 5 is mapped to the lossless `fcoe` forwarding class, and the `fcoe` forwarding class is mapped to the `fabric_fcset_noloss1` fabric fc-set, then the congestion control mechanism is PFC.

- **Case 2**—If the IEEE 802.1p priority 5 is mapped to the lossless fcoe forwarding class, and the fcoe forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop, and the traffic does not receive lossless treatment.
- **Case 3**—If the IEEE 802.1p priority 5 is mapped to the lossless no-loss forwarding class, and the no-loss forwarding class is mapped to the fabric_fcset_noloss2 fabric fc-set, then the congestion control mechanism is PFC.
- **Case 4**—If the IEEE 802.1p priority 5 is mapped to the lossless no-loss forwarding class, and the no-loss forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop, and the traffic does not receive lossless treatment.
- **Case 5**—If the IEEE 802.1p priority 5 is mapped to the lossy best-effort forwarding class, and the best-effort forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop.
- **Case 6**—If the IEEE 802.1p priority 5 is mapped to the lossy best-effort forwarding class, and the best-effort forwarding class is mapped to the fabric_fcset_noloss1 fabric fc-set, then the congestion control mechanism is packet drop.

NOTE: Lossless transport across the fabric must also meet the following two conditions:

1. The maximum cable length between the Node device and the Interconnect device is 150 meters of fiber cable.
2. The maximum frame size is 9216 bytes.

If the MTU is 9216 KB, in some cases the QFabric system supports only five lossless forwarding classes instead of six lossless forwarding classes because of headroom buffer limitations.

The number of IEEE 802.1p priorities (forwarding classes) the QFabric system can support for lossless transport across the Interconnect device fabric depends on several factors:

- **Approximate fiber cable length**—The longer the fiber cable that connects Node device fabric (FTE) ports to the Interconnect device fabric ports, the more data the connected ports need to buffer when a pause is asserted. (The longer the fiber cable, the more frames are traversing the cable when a pause is asserted. Each port must be able to store all of the “in transit” frames in the buffer to preserve lossless behavior and avoid dropping frames.)
- **MTU size**—The larger the maximum frame sizes the buffer must hold, the fewer frames the buffer can hold. The larger the MTU size, the more buffer space each frame consumes.
- **Total number of Node device fabric ports connected to the Interconnect device**—The higher the number of connected fabric ports, the more headroom buffer space the Node device needs on those fabric ports to support the lossless flows that traverse the Interconnect device. Because more buffer

space is used on the Node device fabric ports, less buffer space is available for the Node device access ports, and a lower total number of lossless flows are supported.

The QFabric system supports six lossless priorities (forwarding classes) under most conditions. The priority group headroom that remains after allocating headroom to lossless flows is sufficient to support best-effort and multidestination traffic.

Table 69 on page 326 shows how many lossless priorities the QFabric system supports under different conditions (fiber cable lengths and MTUs) in cases when the QFabric system supports fewer than six lossless priorities. The number of lossless priorities is the same regardless of how many Node device FTE ports are connected to the Interconnect device. However, the higher the number of FTE ports connected to the Interconnect device, the lower the number of total lossless flows supported. In all cases that are not shown in Table 69 on page 326, the QFabric system supports six lossless priorities.

NOTE: The system does not perform a configuration commit check that compares available system resources with the number of lossless forwarding classes configured. If you commit a configuration with more lossless forwarding classes than the system resources can support, frames in lossless forwarding classes might be dropped.

Table 69: Lossless Priority (Forwarding Class) Support for Node Devices When Fewer than Six Lossless Priorities Are Supported

MTU in Bytes	Fiber Cable Length in Meters (Approximate)	Maximum Number of Lossless Priorities (Forwarding Classes) on the Node Device
9216 (9K)	100	5
9216 (9K)	150	5

NOTE: The total number of lossless flows decreases as resource consumption increases. For a Node device, the higher the number of FTE ports connected to the Interconnect device, the larger the MTU, and the longer the fiber cable length, the fewer total lossless flows the QFabric system can support.

RELATED DOCUMENTATION

Understanding CoS Fabric Forwarding Class Sets

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding CoS Scheduling Across the QFabric System | 327](#)

Understanding CoS Scheduling Across the QFabric System

IN THIS SECTION

- [CoS Flow Through the QFabric System | 328](#)
- [Hierarchical Scheduling Architecture on QFabric System Node Devices | 331](#)
- [Default Scheduling on Node Device Fabric Interfaces | 331](#)
- [Hierarchical CoS Architecture Across a QFabric System Interconnect Device | 332](#)
- [Default CoS on Interconnect Device Fabric Interfaces | 334](#)
- [Configuring CoS on Interconnect Device Fabric Interfaces | 342](#)
- [Configuring Scheduling on Node Device Fabric Interfaces | 350](#)
- [Congestion Management | 351](#)

Beginning with Junos OS Release 13.1R2, you can configure two-tier hierarchical scheduling on each Node device fabric interface, and beginning with Junos OS Release 14.1X53-D15, you can configure two-tier hierarchical scheduling on Interconnect device fabric interfaces on a QFabric system. Configuring CoS on the fabric interfaces provides increased control over traffic scheduling across the QFabric system, and helps to ensure predictable bandwidth consumption across the fabric path.

You can configure CoS on the following QFabric system interface types:

- Node device access interfaces (xe interfaces)—Schedule traffic on the output queues of the 10-Gigabit Ethernet access ports using standard Node device CoS scheduling configuration components, as described elsewhere in the QFX Series documentation. You can configure different scheduling for different ports and output queues.
- Node device fabric interfaces (fte interfaces)—Schedule traffic on the output queues of the 40-Gbps fabric interfaces that connect a Node device to a QFX3008-I or a QFX3600-I Interconnect device using standard Node device CoS scheduling configuration components. You can configure different scheduling for different interfaces and output queues.
- Interconnect device fabric interfaces (fte interfaces)—Schedule traffic on the output queues of the 40-Gbps fabric interfaces that connect an Interconnect device to a Node device. Configuring

schedulers, mapping schedulers to output queues, and applying scheduling to interfaces on Interconnect devices differ in some aspects from scheduling configuration on Node devices. You can configure different scheduling for different interfaces and fabric forwarding class sets (fabric fc-sets).

- Interconnect device internal Clos fabric interfaces (bfte interfaces)—Schedule traffic on the internal 40-Gbps Clos fabric interfaces that connect the ingress and egress stages of the Interconnect device Clos fabric, using the same scheduling components as the Interconnect device fabric (fte) interfaces. You can configure one Clos fabric interface scheduler, which the system applies to all of the internal Clos fabric interfaces. You cannot configure different schedulers for different Clos fabric interfaces.

Configuring scheduling on Interconnect device fabric interfaces differs from configuring scheduling on Node device interfaces because the Interconnect device is a shared infrastructure that supports traffic from multiple Node devices and CoS configurations.

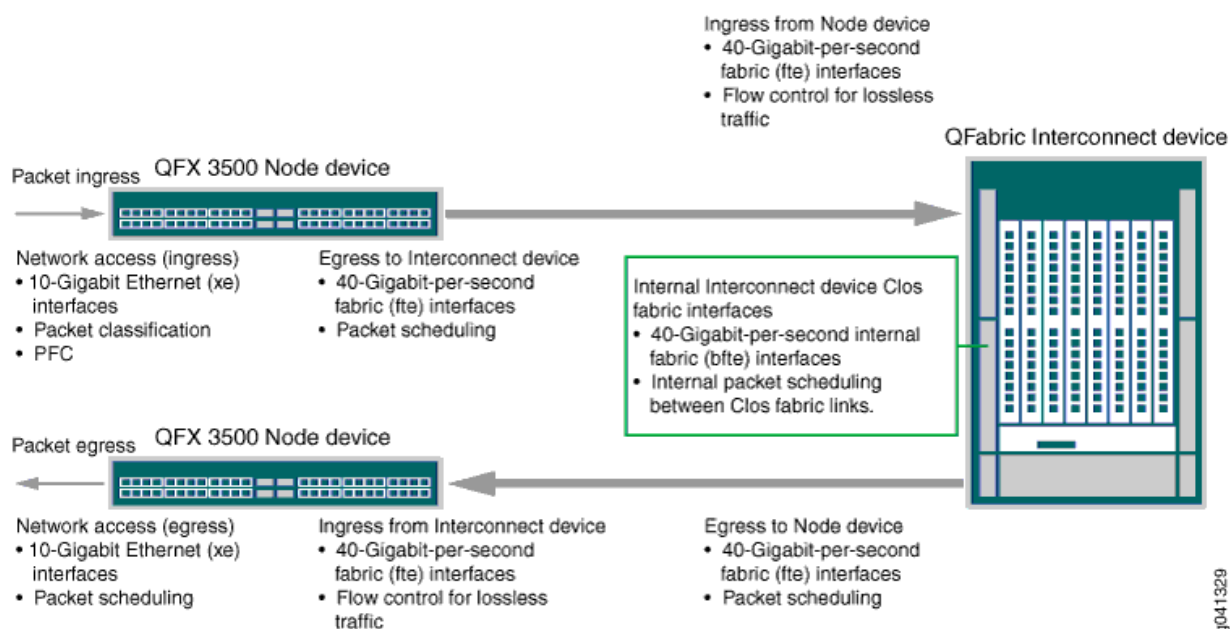
NOTE: On Node device access interfaces and fabric, the hierarchical scheduling you configure is the Junos OS implementation of enhanced transmission selection (ETS, described in IEEE 802.1Qaz). On Interconnect device fabric interfaces, the hierarchical scheduling you configure is not an implementation of ETS, although it functions similarly to ETS in that excess port bandwidth is shared.

If the 40-Gbps fabric links that connect Node devices to Interconnect devices become oversubscribed, you can configure CoS to control how those fabric links allocate bandwidth to traffic as described in this topic:

CoS Flow Through the QFabric System

[Figure 13 on page 329](#) shows the CoS flow across the QFabric system.

Figure 13: QFabric System CoS Flow



Packets from the access network enter the QFabric system at the ingress interfaces of a QFabric system Node device, cross the Interconnect device fabric, and then are forwarded to their destination through the egress interfaces of another QFabric system Node device.

NOTE: Traffic that uses the same Node device for both traffic ingress and traffic egress does not cross the fabric. CoS for this type of traffic is the same as CoS on a standalone switch.

When a packet enters the QFabric system, it receives CoS treatment at each interface it traverses:

1. A packet enters the QFabric system on a 10-Gigabit Ethernet access interface on a QFabric Node device. At the Node device ingress interface, the packet is classified into a forwarding class, which groups the packet with other traffic that requires similar CoS treatment and maps the packet to the appropriate output queue. To support lossless traffic delivery, enable PFC on the IEEE 802.1p code points of lossless priorities.
2. Next, the packet exits the QFabric Node device on a 40-Gbps fabric interface that is connected to the QFabric Interconnect device. At the Node device egress interface, the packet is placed in the correct output queue and receives the configured (or default) CoS scheduling, which determines the bandwidth and priority allocated to the packet for its journey from the Node device to the Interconnect device.
3. The packet enters the Interconnect device on the 40-Gbps fabric interface connected to the ingress Node device. At the Interconnect device ingress interface, the forwarding class of the packet maps the packet to a fabric fc-set, which groups the packet with other traffic that requires similar CoS

treatment and maps the packet to the appropriate output queue. Flow control is applied automatically to traffic in lossless fabric fc-sets to preserve the lossless characteristics of that traffic. (Lossless forwarding classes are mapped to the lossless fabric fc-sets.) Other traffic uses standard packet drop for flow control.

4. The packet progresses from the Interconnect device ingress interface to the internal, three-stage, 40-Gbps Clos fabric interfaces. At the Clos fabric interfaces, packet flow control to protect lossless traffic is applied automatically to traffic in lossless fabric fc-sets. At the egress interfaces from the Clos fabric interfaces, the packet is placed in the correct output queue and receives CoS scheduling.

NOTE: If you do not use the default Clos fabric interface scheduling, you can configure one scheduler that is applied to all three of the Clos fabric interfaces.

5. The packet exits the QFabric Interconnect device on the 40-Gbps fabric interface connected to the egress Node device. At the Interconnect device egress interface, the packet is placed in the correct output queue and receives the configured (or default) CoS scheduling, which determines the bandwidth and priority allocated to the packet for its journey from the Interconnect device egress to the Node device.
6. The packet enters the egress Node device on the 40-Gbps fabric interface connected to the Interconnect device egress interface. At the Node device fabric interface, the packet forwarding class determines the fc-set in which the packet is placed and the output queue the packet uses. Packet flow control to protect lossless traffic is applied automatically to traffic in lossless fabric fc-sets.
7. The packet exits the QFabric system from the egress Node device on a 10-Gigabit Ethernet access interface. At the Node device egress interface, the packet is placed in the correct output queue and receives the configured (or default) CoS scheduling, which determines the bandwidth and priority allocated to the packet for its journey from the Node device to the packet destination.

You can use default CoS scheduling or configure CoS scheduling on any or all of the Node device interfaces and on Interconnect device fabric (*fte*) interfaces. If you configure scheduling on one of these interfaces, you can still use default scheduling on other interfaces. Because you configure one scheduler for all of the Interconnect device Clos fabric interfaces (*bfte* interfaces), all of the Clos fabric interfaces either use the configured scheduling or the default scheduling, but not a mix of configured and default scheduling.

NOTE: To support lossless traffic delivery, you must enable PFC on the IEEE 802.1p code points of lossless priorities (forwarding classes) at the Node device network access ingress interfaces. Flow control is applied to lossless priorities automatically on the fabric (*fte* and *bfte*) interfaces.

Hierarchical Scheduling Architecture on QFabric System Node Devices

CoS architecture on Node device access interfaces is the same as CoS architecture on standalone switch access interfaces. CoS architecture on Node device fabric interfaces is also the same as the CoS architecture on the access interfaces. You apply schedulers to queues (priorities), fc-sets (priority groups), and interfaces in the same hierarchical manner as described in ["Understanding CoS Hierarchical Port Scheduling \(ETS\)" on page 223](#).

You configure scheduling on Node device fabric interfaces (fte interfaces) using the same statements and configuration constructs that you use to configure scheduling on Node device access interfaces (xe interfaces). For example, on Node device fabric interfaces you can:

- Define up to four fc-sets (three unicast, one multidestination)

NOTE: If the interface handles strict-high priority traffic, you must define a separate fc-set (priority group) for strict-high priority traffic. Strict-high priority traffic cannot be mixed with traffic of other priorities in an fc-set. For example, you might choose to create different fc-sets for best effort, lossless, strict-high priority, and multidestination traffic.

- Map forwarding classes to fc-sets
- Configure scheduling for each forwarding class (scheduler)
- Configure scheduling for each fc-set (traffic control profile)

The only differences in configuring CoS on Node device fabric interfaces compared to configuring CoS on Node device access interfaces are:

- You specify a Node device fabric interface instead of a Node device access interface when you apply CoS to an interface.
- You cannot attach classifiers or congestion notification profiles to fabric interfaces.

Default Scheduling on Node Device Fabric Interfaces

Default scheduling on Node device fabric interfaces is the same as default scheduling on Node device access interfaces. Only the default forwarding classes (best-effort, network-control, fcoe, no-loss, and multidestination) receive port bandwidth, based on the default minimum guaranteed bandwidth (transmit rate) scheduler settings for each default forwarding class.

All of the default forwarding classes are placed in one default group and receive port bandwidth based on their default transmit rate settings (weights). Forwarding classes that are not default forwarding classes receive no bandwidth.

Each default forwarding class receives a guaranteed minimum percentage of the port bandwidth based on the default transmit rate. [Table 70 on page 332](#) shows the default transmit rate for each of the default forwarding classes.

Table 70: Default Node Device Fabric Interface Forwarding Class Scheduler Configuration

Default Forwarding Classes	Transmit Rate (Percentage of Class Group Bandwidth)
best-effort	5%
fcoe	35%
no-loss	35%
network-control	5%
mcast	20%

Bandwidth is divided among the default forwarding classes in a ratio proportional to the default transmit rate for the forwarding class.

Hierarchical CoS Architecture Across a QFabric System Interconnect Device

Because Interconnect devices support traffic from multiple Node devices that have multiple CoS configurations, configuring CoS on Interconnect device fabric interfaces differs from configuring CoS on Node device access and fabric interfaces.

The hierarchical CoS scheduling structure on the Interconnect device interfaces consists of two tiers:

1. Fabric forwarding class sets—Similar to fc-sets on Node devices, fabric fc-sets group traffic for transport across the Interconnect device fabric. Fabric fc-sets are global and apply to all traffic that crosses the fabric from all Node devices. See ["Understanding CoS Fabric Forwarding Class Sets" on page 396](#) for a detailed description of fabric fc-sets.
2. Class groups—Fabric fc-sets are grouped into class groups for transport across the Interconnect device.

Node devices and Interconnect devices each have a two-tier hierarchical CoS scheduling architecture. The architectures are slightly different, but each scheduling tier performs analogous functions, as shown in [Table 71 on page 333](#).

Table 71: Bandwidth Scheduler Architecture on Node Devices and Interconnect Devices

Bandwidth Allocation Pool	Node Devices	Interconnect Devices
Port—Entire amount of bandwidth available to traffic on a port.	Access (xe) or fabric (fte) interfaces	Fabric (fte) or Clos fabric (bfte) interfaces
Priority group—Group of traffic types that requires similar CoS treatment. Each priority group receives a portion of the total available port bandwidth.	Forwarding class set (fc-set)	Class group
Priority—Most granular tier of bandwidth allocation. Each priority receives a portion of the total available priority group bandwidth.	Forwarding class (mapped to output queue)	Fabric fc-set (mapped to output queue)

Fabric FC-Sets

Fabric fc-sets are groups of forwarding classes that receive similar CoS treatment across the Interconnect device. Fabric fc-sets are global to the QFabric system and apply to all traffic that traverses the fabric, from all connected Node devices. The CoS you configure on a fabric fc-set applies to all the traffic that belongs to that fabric fc-set.

For example, a fabric fc-set that includes the best-effort forwarding class handles all of the best-effort traffic from all of the connected Node devices that traverses the Interconnect device fabric.

There are 12 default fabric fc-sets, including five visible fabric fc-sets and seven hidden fabric fc-sets. The five visible fabric fc-sets have forwarding classes mapped to them by default. By default, the seven hidden fabric fc-sets do not carry traffic, but you can map forwarding classes to the hidden fabric fc-sets if you want to use them.

You can configure the forwarding class membership of each fabric fc-set, and you can configure CoS for each fabric fc-set. However, you cannot create new fabric fc-sets, and you cannot delete the 12 default fabric fc-sets.

Each fabric fc-set is mapped to an output queue. Each fabric interface has 12 output queues, one for each of the 12 fabric fc-sets. The traffic from all of the forwarding classes mapped to a fabric fc-set uses that fabric fc-set's output queue.

Fabric fc-sets are grouped into class groups for transport across the Interconnect device.

Class Groups for Fabric FC-Sets

To transport traffic across the fabric, the QFabric system organizes the fabric fc-sets into three default classes called *class groups*. Class groups are not user-configurable. The three class groups are:

- **Strict-high priority**—All traffic in the fabric fc-set `fabric_fcset_strict_high`. This class group includes the traffic in strict-high priority and network-control forwarding classes and in any forwarding classes you create on a Node device that consist of strict-high priority or network-control forwarding class traffic.
- **Unicast**—All traffic in the fabric fc-sets `fabric_fcset_be`, `fabric_fcset_noloss1`, and `fabric_fcset_noloss2`. This class group includes the traffic in the best-effort, fcoe, and no-loss forwarding classes, and the traffic in any forwarding classes you create on a Node device that consist of best-effort or lossless traffic. If you use any of the hidden no loss fabric fc-sets (`fabric_fcset_noloss3`, `fabric_fcset_noloss4`, `fabric_fcset_noloss5`, or `fabric_fcset_noloss6`), that traffic is part of this class group.
- **Multidestination**—All traffic in the fabric fc-set `fabric_fcset_multicast1`. This class group includes the traffic in the mcast forwarding class and in any forwarding classes you create on a Node device that consist of multidestination traffic. If you use any of the hidden multidestination fabric fc-sets (`fabric_fcset_multicast2`, `fabric_fcset_multicast3`, or `fabric_fcset_multicast4`), that traffic is part of this class group.

Default CoS on Interconnect Device Fabric Interfaces

If you do not configure CoS on the Interconnect device fabric interfaces, the Interconnect device interfaces use the default CoS configuration as described in this section:

Default Class Group Scheduling

Default class group bandwidth scheduling is analogous to default fc-set (priority group) scheduling on a Node device. Default class group scheduling uses weighted round-robin (WRR) scheduling, in which each class group receives a portion of the total available fabric interface bandwidth, based on the class group's traffic type, as shown in [Table 72 on page 335](#):

Table 72: Class Group Default Scheduling Properties and Membership

Class Group	Fabric fc-sets	Forwarding Classes (Default Mapping)	Class Group Scheduling Properties (Weight)
Strict-high priority	fabric_fcset_strict_high	<ul style="list-style-type: none"> All strict-high priority forwarding classes network-control 	Traffic in the strict-high priority class group is served first. This class group receives all of the bandwidth it needs to empty its queues and therefore can starve other types of traffic during periods of high-volume strict priority traffic. Plan carefully and use caution when determining how much traffic to configure as strict-high priority traffic.
Unicast	<ul style="list-style-type: none"> fabric_fcset_be fabric_fcset_nolos1 fabric_fcset_nolos2 Includes the hidden lossless fabric fc-sets if used: <ul style="list-style-type: none"> fabric_fcset_nolos3 fabric_fcset_nolos4 fabric_fcset_nolos5 fabric_fcset_nolos6 	<ul style="list-style-type: none"> best-effort fcoe no-loss NOTE: No forwarding classes are mapped to the hidden lossless fabric_fcsets by default.	Traffic in the unicast class group receives an 80% weight in the weighted round-robin (WRR) calculations. After the strict-high priority class group has been served, the unicast class group receives 80% of the remaining fabric bandwidth. (If more bandwidth is available, the unicast class group can use more bandwidth.)

Table 72: Class Group Default Scheduling Properties and Membership (Continued)

Class Group	Fabric fc-sets	Forwarding Classes (Default Mapping)	Class Group Scheduling Properties (Weight)
Multidestination	<p>fabric_fcset_multicast1</p> <p>Includes the hidden multidestination fabric fc-sets if used:</p> <ul style="list-style-type: none"> • fabric_fcset_multicast2 • fabric_fcset_multicast3 • fabric_fcset_multicast4 	<ul style="list-style-type: none"> • mcast <p>NOTE: No forwarding classes are mapped to the hidden multidestination fabric_fcsets by default.</p>	Traffic in the multidestination class group receives a 20% weight in the WRR calculations. After the strict-high priority class group has been served, the multidestination class group receives 20% of the remaining fabric bandwidth. (If more bandwidth is available, the multidestination class group can use more bandwidth.)

If you use the default fabric CoS configuration, only the five visible fabric fc-sets have traffic mapped to them by default. The fabric fc-sets within each class group are weighted by their transmit rates (guaranteed minimum bandwidth), and they receive bandwidth from the class group's total bandwidth using weighted round-robin (WRR) scheduling.

Default Fabric FC-Set Bandwidth Scheduling

Default fabric fc-set bandwidth scheduling is analogous to default forwarding class (priority) scheduling on a Node device. Each fabric fc-set receives a guaranteed minimum percentage of the port bandwidth that the class group receives. The guaranteed minimum percentage is called the *transmit rate*.

[Table 73 on page 336](#) shows the default transmit rate for each of the default fabric fc-sets.

Table 73: Default Fabric FC-Set Scheduler Configuration

Default Fabric FC-Set	Transmit Rate (Percentage of Class Group Bandwidth)
fabric_fcset_strict_high	0%
fabric_fcset_noloss1	35%

Table 73: Default Fabric FC-Set Scheduler Configuration (Continued)

Default Fabric FC-Set	Transmit Rate (Percentage of Class Group Bandwidth)
fabric_fcset_noloss2	35%
fabric_fcset_be	10%
fabric_fcset_multicast1	20%

Each fabric fc-set belongs to a class group. Each class group receives a portion of the total available port bandwidth. Each fabric fc-set in a class group receives a portion of the total available class group bandwidth based on the transmit rate (weight) of the fabric fc-set.

Traffic in fabric_fcset_strict_high does not have a default transmit rate because fabric_fcset_strict_high receives all of the bandwidth needed to empty its queue before other queues are served. Traffic in the remaining fabric fc-sets receive bandwidth in a ratio proportional to the default transmit rate of each fabric fc-set.

Each of the following hidden fabric fc-sets receives a default scheduling weight of 1 if you do not configure CoS scheduling for it:

- fabric_fcset_noloss3
- fabric_fcset_noloss4
- fabric_fcset_noloss5
- fabric_fcset_noloss6
- fabric_fcset_multicast2
- fabric_fcset_multicast3
- fabric_fcset_multicast4

You must explicitly map forwarding classes to hidden fabric fc-sets and configure scheduling for that traffic if you want to use the hidden fabric fc-sets. Default scheduling does not use the hidden fabric fc-sets.

Default Class Group and Fabric FC-Set Scheduling Example

The following example shows how default scheduling allocates the total port bandwidth among the class groups and fabric fc-sets. The example assumes that traffic is mapped to each of the forwarding classes

in the five visible fabric fc-sets, and that the strict-high priority class group consumes an average of 10 percent of the 40-Gbps fabric interface bandwidth (4 gigabits), leaving 90 percent of the fabric interface bandwidth (36 gigabits) for the remaining class groups.

In this scenario, by default, the strict-high priority class group includes one fabric fc-set (fabric_fcset_strict_high), the unicast class group includes three fabric fc-sets (fabric_fcset_be, fabric_fcset_noloss1, and fabric_fcset_noloss2) and the multideestination class group includes one fabric fc-set (fabric_fcset_multicast1). Each individual fabric fc-set receives the following treatment:

- Strict-high priority class group (fabric_fcset_strict_high)—Assumed to average 10 percent (4 gigabits) for the purposes of this example. Because the strict-high priority class group is served first and receives all of the bandwidth it requires to empty its queue, in real networks the amount of required bandwidth fluctuates and affects the amount of bandwidth available to the other class groups.

TIP: To prevent strict-high priority traffic from using too much bandwidth, you can set a maximum bandwidth limit by configuring a scheduler shaping rate for the fabric_fcset_strict_high fabric fc-set.

- Unicast class group (fabric_fcset_be, fabric_fcset_noloss1, and fabric_fcset_noloss2)—Each of these fabric fc-sets receives a weighted portion of the 80 percent of the total port bandwidth available to the unicast class group after the strict-high traffic has been served. The weight corresponds to the transmit rate of each fabric fc-set. The following calculations show the minimum port bandwidth allocated to each of the unicast class group fabric fc-sets:

- fabric_fcset_be

$10/(35+35+10)\%$ of 80% of the available port bandwidth (12.5% of 80% of port bandwidth)

The 10 that is the numerator in $10/(35+35+10)$ is the percentage of bandwidth allocated to the fabric_fcset_be by the transmit rate weight. The $(35+35+10)$ in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the three fabric fc-sets in the unicast class group.

The 80 percent represents 80 percent of the port bandwidth available after serving strict-high priority traffic (36 gigabits).

The resulting equation is:

$10/(35+35+10)\% \times (0.8 \times 36 \text{ gigabits}) = \text{approximately } 3.6 \text{ gigabits}$

- fabric_fcset_noloss1 and fabric_fcset_noloss2

The default minimum bandwidth for the two visible lossless fabric fc-sets is the same because both of these fabric fc-sets have the same transmit rate weight.

$35/(35+35+10)\%$ of 80% of the port bandwidth (43.75% of 80% of port bandwidth)

The 35 that is the numerator in $35/(35+35+10)$ is the percentage of bandwidth allocated to each of the no-loss fabric fc-sets by the transmit rate weight. The $(35+35+10)$ in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the three fabric fc-sets in the unicast class group.

The 80 percent represents 80 percent of the port bandwidth available after serving strict-high priority traffic (36 gigabits).

The resulting equation is:

$$35/(35+35+10)\% \times (0.8 \times 36 \text{ gigabits}) = \text{approximately } 12.6 \text{ gigabits}$$

- Multidestination class group (fabric_fcset_multicast1)—Because only one fabric fc-set is configured by default in the multidestination class group, it receives 100 percent of the 20 percent of the total port bandwidth available to the multidestination class group after the strict-high traffic has been served:

$$100/(100)\% \text{ of } 20\% \text{ of the available port bandwidth (100\% of 20\% of available port bandwidth)}$$

The resulting equation is:

$$100/100\% \times (0.2 \times 36 \text{ gigabits}) = \text{approximately } 7.2 \text{ gigabits}$$

Default PFC and Lossless Transport Across the Interconnect Device

The Interconnect device incorporates flow control mechanisms to support lossless transport during periods of congestion on the fabric. To support the priority-based flow control (PFC) feature on the Node devices, the Interconnect device fabric supports lossless transport for up to six IEEE 802.1p priorities when the following two configuration constraints are met:

1. The IEEE 802.1p priority used for the traffic that requires lossless transport is mapped to a lossless forwarding class (a forwarding class configured with the `no-loss` parameter or the default `fcoe` or `no-loss` forwarding class).
2. The lossless forwarding class must be mapped to one of the lossless fabric fc-sets (fabric_fcset_noloss1, fabric_fcset_noloss2, fabric_fcset_noloss3, fabric_fcset_noloss4, fabric_fcset_noloss5, or fabric_fcset_noloss6). If you do not explicitly map lossless forwarding classes to fabric fc-sets, lossless forwarding classes are mapped by default to lossless fabric fc-sets fabric_fcset_noloss1 and fabric_fcset_noloss2.

When traffic meets these two constraints, the fabric propagates back-pressure from egress queues during periods of congestion. However, to achieve end-to-end lossless transport across the QFabric system, you must also configure a congestion notification profile to enable PFC on the Node device ingress interfaces. To achieve end-to-end lossless transport across the network, you must configure PFC on all of the devices in the lossless traffic path.

For all other combinations of IEEE 802.1p priority to forwarding class mapping and all other combinations of forwarding class to fabric fc-set mapping, the default congestion control mechanism is normal packet drop. For example:

- **Case 1**—If the IEEE 802.1p priority 5 is mapped to the lossless fcoe forwarding class, and the fcoe forwarding class is mapped to the fabric_fcset_noloss1 fabric fc-set, then the congestion control mechanism is PFC.
- **Case 2**—If the IEEE 802.1p priority 5 is mapped to the lossless fcoe forwarding class, and the fcoe forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop, and the traffic does not receive lossless treatment.
- **Case 3**—If the IEEE 802.1p priority 5 is mapped to the lossless no-loss forwarding class, and the no-loss forwarding class is mapped to the fabric_fcset_noloss2 fabric fc-set, then the congestion control mechanism is PFC.
- **Case 4**—If the IEEE 802.1p priority 5 is mapped to the lossless no-loss forwarding class, and the no-loss forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop, and the traffic does not receive lossless treatment.
- **Case 5**—If the IEEE 802.1p priority 5 is mapped to the lossy best-effort forwarding class, and the best-effort forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop.
- **Case 6**—If the IEEE 802.1p priority 5 is mapped to the lossy best-effort forwarding class, and the best-effort forwarding class is mapped to the fabric_fcset_noloss1 fabric fc-set, then the congestion control mechanism is packet drop.

NOTE: Lossless transport across the fabric must also meet the following two conditions:

1. The maximum cable length between the Node device and the Interconnect device is 150 meters of fiber cable.
2. The maximum frame size is 9216 bytes.

If the MTU is 9216 KB, in some cases the QFabric system supports only five lossless forwarding classes instead of six lossless forwarding classes because of headroom buffer limitations.

The number of IEEE 802.1p priorities (forwarding classes) the QFabric system can support for lossless transport across the Interconnect device fabric depends on several factors:

- **Approximate fiber cable length**—The longer the fiber cable that connects Node device fabric (FTE) ports to the Interconnect device fabric ports, the more data the connected ports need to buffer when a pause is asserted. (The longer the fiber cable, the more frames are traversing the cable when

a pause is asserted. Each port must be able to store all of the “in transit” frames in the buffer to preserve lossless behavior and avoid dropping frames.)

- MTU size—The larger the maximum frame sizes the buffer must hold, the fewer frames the buffer can hold. The larger the MTU size, the more buffer space each frame consumes.
- Total number of Node device fabric ports connected to the Interconnect device—The higher the number of connected fabric ports, the more headroom buffer space the Node device needs on those fabric ports to support the lossless flows that traverse the Interconnect device. Because more buffer space is used on the Node device fabric ports, less buffer space is available for the Node device access ports, and a lower total number of lossless flows are supported.

The QFabric system supports six lossless priorities (forwarding classes) under most conditions. The priority group headroom that remains after allocating headroom to lossless flows is sufficient to support best-effort and multidestination traffic.

[Table 74 on page 341](#) shows how many lossless priorities the QFabric system supports under different conditions (fiber cable lengths and MTUs) in cases when the QFabric system supports fewer than six lossless priorities. The number of lossless priorities is the same regardless of how many Node device FTE ports are connected to the Interconnect device. However, the higher the number of FTE ports connected to the Interconnect device, the lower the number of total lossless flows supported. In all cases that are not shown in [Table 74 on page 341](#), the QFabric system supports six lossless priorities.

NOTE: The system does not perform a configuration commit check that compares available system resources with the number of lossless forwarding classes configured. If you commit a configuration with more lossless forwarding classes than the system resources can support, frames in lossless forwarding classes might be dropped.

Table 74: Lossless Priority (Forwarding Class) Support for Node Devices When Fewer than Six Lossless Priorities Are Supported

MTU in Bytes	Fiber Cable Length in Meters (Approximate)	Maximum Number of Lossless Priorities (Forwarding Classes) on the Node Device
9216 (9K)	100	5
9216 (9K)	150	5

NOTE: The total number of lossless flows decreases as resource consumption increases. For a Node device, the higher the number of FTE ports connected to the Interconnect device, the larger the MTU, and the longer the fiber cable length, the fewer total lossless flows the QFabric system can support.

Configuring CoS on Interconnect Device Fabric Interfaces

If you do not want to use default CoS scheduling across the Interconnect device fabric, you can configure two-tier hierarchical scheduling on the external 40-Gbps fabric interfaces (fte interfaces) and on the internal 40-Gbps Clos fabric interfaces (bfte interfaces).

This section describes:

Similarities Between Node Device Scheduling and Interconnect Device Scheduling

Configuring two-tier hierarchical scheduling on Interconnect device fabric interfaces follows the same general process as configuring scheduling on Node device interfaces, in that you perform the following actions in both cases:

- Define drop profiles to control packet loss for lossy traffic; do not use drop profiles on lossless traffic or multideestination traffic. (However, if you configure a drop profile on lossless traffic or on multideestination traffic, the system does not return a commit error.)
- Define schedulers to configure the bandwidth for different types of traffic.
- Map schedulers to output queues (by mapping schedulers to forwarding classes on Node devices, and by mapping schedulers to fabric fc-sets on Interconnect devices).
- Associate hierarchical scheduling with interfaces to apply scheduling to traffic on those interfaces.

Another similarity is that you cannot configure classifiers or congestion notification profiles (to enable PFC) on fabric interfaces. Flow control is applied automatically to lossless queues on fabric interfaces, and packet classification occurs at the Node device ingress interface.

Differences Between Node Device and Interconnect Device Hierarchical Scheduling

Configuring the two-tier scheduling hierarchy on Interconnect device fabric interfaces is different in several important ways than configuring the two-tier scheduling hierarchy on Node device interfaces, as shown in [Table 75 on page 343](#):

Table 75: Node Device and Interconnect Device Hierarchical Scheduling Differences

Hierarchical Scheduling Component	Node Devices	Interconnect Devices
Priority scheduling hierarchy tier	<ul style="list-style-type: none"> Each forwarding class is mapped to an output queue. Classifiers map forwarding classes to priorities (IEEE 802.1p code points). You map schedulers to forwarding classes to provide scheduling for priorities. You can create and delete forwarding classes. 	<ul style="list-style-type: none"> Each fabric fc-set is mapped to an output queue, and is mapped internally to priorities (IEEE 802.1p code points). You map schedulers to fabric fc-sets to provide scheduling for priorities. You cannot create or delete fabric fc-sets. Only the 12 default fabric fc-sets are available (but you can change the default mapping of forwarding classes to fabric fc-sets).
Scheduler mapping to priorities (bandwidth allocation to priorities)	The scheduler-maps statement maps a forwarding class to a scheduler.	The scheduler-map-forwarding-class-sets statement maps a fabric fc-set to a scheduler.
Priority group scheduling hierarchy tier	<ul style="list-style-type: none"> Each fc-set represents a priority group. You associate fc-sets with traffic control profiles to provide scheduling for priority groups. You can create and delete fc-sets. 	<ul style="list-style-type: none"> Each class group represents a priority group. You cannot change the types of traffic associated with a class set (each class set is dedicated to one type of traffic: strict-high priority, unicast, or multidestination traffic). You cannot create or delete class groups.

Table 75: Node Device and Interconnect Device Hierarchical Scheduling Differences *(Continued)*

Hierarchical Scheduling Component	Node Devices	Interconnect Devices
Priority group bandwidth allocation method	You create traffic control profiles to determine the port scheduling resources assigned to priority groups (fc-sets).	You do not configure priority group (class group) scheduling using a traffic control profile. Instead, the QFabric system uses the sum of the fabric fc-set minimum guaranteed bandwidths (transmit rates) to determine the port scheduling resources for the class group, as described in "Hierarchical Scheduling Bandwidth Allocation" on page 346 later in this topic.
Scheduler transmit rate, shaping rate, and drop priority parameters	<ul style="list-style-type: none"> • Transmit rate and shaping rate—You can specify either a percentage value or an absolute value for these two parameters. • Priority—Scheduling for forwarding classes includes the priority parameter, which sets the scheduling drop priority as either low or strict-high. 	<ul style="list-style-type: none"> • Transmit rate and shaping rate—You can only specify a percentage value for these two parameters; you cannot specify an absolute value. • Priority—You cannot specify the priority parameter because the class groups automatically determine the drop priority. If you try to map a scheduler that includes a priority setting to a fabric fc-set, the system generates a commit error.
Hierarchical scheduler association with interfaces	Specify an access interface or a Node device fabric interface and associate it with an fc-set (determines which forwarding classes use the interface) and a traffic control profile (determines scheduling for both the priority group and the priorities in the priority group).	<p>Specify an Interconnect device fabric interface and associate it with a fabric forwarding class set scheduler map. The fabric forwarding class set scheduler map determines the fabric fc-sets associated with the interface and their scheduling properties.</p> <p>You can associate one fabric forwarding class scheduler map with an interface. Different interfaces can have different fabric forwarding class scheduler maps.</p>

Table 75: Node Device and Interconnect Device Hierarchical Scheduling Differences *(Continued)*

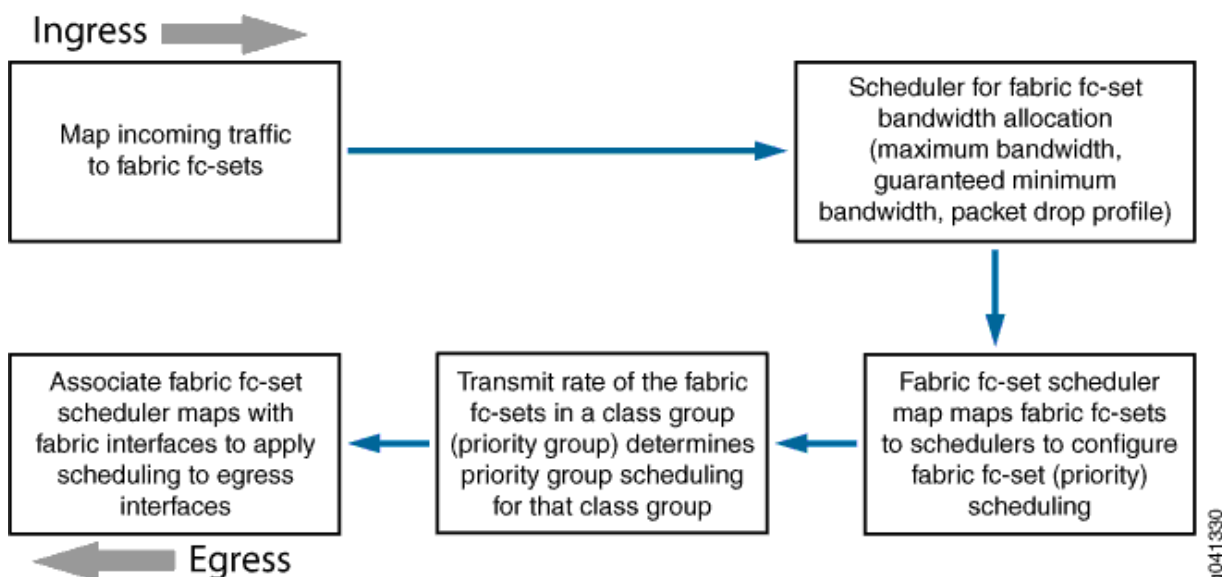
Hierarchical Scheduling Component	Node Devices	Interconnect Devices
Classifiers and congestion notification profiles (enabling PFC on priorities)	You can attach classifiers and congestion notification profiles to access interfaces (although you cannot attach them to fabric interfaces).	You cannot attach classifiers and congestion notification profiles to fabric interfaces.

NOTE: Because the queue scheduler transmit rate is used differently on Node devices and Interconnect devices, and because you cannot specify the scheduler priority parameter on Interconnect devices, you should configure different schedulers for Node device interfaces and Interconnect device interfaces.

Hierarchical Scheduling Configuration Components

Some of the configuration components used for Interconnect device CoS scheduling are similar to the CoS configuration components used for Node device CoS scheduling, but some of the components are different because configuring the two-tier scheduling hierarchy differs in some respects on the two devices. [Figure 14 on page 346](#) shows a block diagram of the components used to configure hierarchical scheduling on the Interconnect device.

Figure 14: Configuration Components of Interconnect Device Hierarchical Scheduling



Hierarchical Scheduling Bandwidth Allocation

The purpose of hierarchical scheduling is to allocate the available port bandwidth to class groups (priority groups), and then to allocate class group bandwidth to the fabric fc-sets (priorities) that belong to the class group. Hierarchical scheduling provides better port bandwidth utilization and greater flexibility to allocate port resources to queues (priorities) and to groups of queues (priority groups) than flat scheduling.

NOTE: Available port bandwidth is the bandwidth that remains after the port services all of its strict-high priority traffic.

You allocate bandwidth to priorities by configuring scheduling for fabric fc-sets. For each fabric fc-set, you can configure a scheduler that defines the guaranteed minimum bandwidth (transmit rate), the maximum bandwidth (shaping rate), and the packet drop profile for lossy unicast traffic assigned to that fabric fc-set. (Lossless fabric fc-sets use flow control to prevent packet loss and do not use drop profiles; multideestination traffic does not use drop profiles.)

Bandwidth is allocated to priority groups (class groups) automatically, based on the minimum guaranteed bandwidth (transmit rate) of the fabric fc-sets that belong to the class group. The sum of the transmit rates of the fabric fc-sets in a class group equals the total minimum guaranteed port bandwidth of that class group.

So the QFabric system uses the fabric fc-set transmit in two ways to calculate bandwidth allocation:

1. The transmit rate of a fabric fc-set sets the minimum guaranteed bandwidth allocated to that fabric fc-set from the class group bandwidth pool.
2. The sum of the fabric fc-set transmit rates in a class group sets the minimum guaranteed port bandwidth allocated to that class group.

The transmit rate percentage that you configure in a fabric fc-set scheduler does not necessarily equal the minimum percentage of available port bandwidth allocated to that fabric fc-set, because port bandwidth is allocated to strict-high priority traffic first, and only the remaining port bandwidth is allocated to the rest of the traffic based on the fabric fc-set transmit rates. In other words, the bandwidth available to a class group after the system services strict-high priority traffic is divided among the fabric fc-sets in that class group in proportion to the transmit rate configured for each fabric fc-set.

Hierarchical scheduling on fabric interfaces allocates guaranteed minimum port bandwidth in the following manner:

1. The sum of the transmit rates of the fabric fc-sets in a class group determines the amount of available port bandwidth allocated to the class group. For example, a class group that has three fabric fc-sets with transmit rates of 10 percent, 20 percent, and 30 percent, receives 60 percent of the available port bandwidth ($10+20+30 = 60$).
2. The fabric fc-set transmit rate is used again to determine the proportion of class group bandwidth allocated to the fabric fc-set. For example, in a class group with three fabric fc-sets that have transmit rates of 10 percent, 20 percent, and 30 percent (class group receives 60 percent of available port bandwidth), the fabric fc-set with a transmit rate of 20 percent receives one-third of the class group bandwidth (20 is one-third of 60).

It is important to understand that this is not one-third of the total available port bandwidth, but one-third of the 60 percent of total available port bandwidth that the class group receives.

NOTE: The sum of the transmit rates of all of the fabric fc-sets in the unicast and the multidestination class groups cannot exceed 100 percent. (You cannot configure the system to schedule more than 100 percent as the minimum guaranteed bandwidth for all of the unicast and multidestination fabric fc-sets. The sum of the transmit rates of all unicast and multidestination fabric fc-sets must be less than or equal to 100 percent.)

Interconnect Device Hierarchical Scheduling (Class Group and Fabric FC-Set) Example

The following example shows how configuring hierarchical scheduling allocates the total port bandwidth among the class groups and fabric fc-sets. The example shows a configuration in which:

- The strict-high priority class group has no scheduler (transmit rate of fabric_fcset_strict_high is 0 percent and no maximum bandwidth is set, so the strict-high priority traffic can use as much bandwidth as needed).
- The unicast class group consists of the following three fabric fc-sets and transmit rates:
 - fabric_fcset_be, 25 percent
 - fabric_fcset_noloss1, 15 percent
 - fabric_fcset_noloss2, 20 percent
- The multideestination class group consists of the following two fabric fc-sets and transmit rates:
 - fabric_fcset_multicast1, 10 percent
 - fabric_fcset_multicast2, 30 percent

Total available port bandwidth (port bandwidth remaining after serving strict-high priority traffic) is divided between the unicast and multideestination class groups:

- Unicast class group—Contains three fabric fc-sets (fabric_fcset_be, fabric_fcset_noloss1, and fabric_fcset_noloss2) with a combined transmit rate of 60 percent (25+15+ 20). Therefore, the unicast class group receives 60 percent of the total available port bandwidth.
- Multideestination class group—Contains two fabric fc-sets (fabric_fcset_multicast1 and fabric_fcset_multicast2) with a combined transmit rate of 40 percent (10+30). Therefore, the multideestination class group receives 40 percent of the total available port bandwidth.

The class group bandwidth is divided among the fabric fc-sets based on the transmit rate of each fabric fc-set in relation to the class group bandwidth.

The unicast class group bandwidth is divided among its three fabric fc-sets:

- fabric_fcset_be

25/(15+20+25) percent of 60 percent of the available port bandwidth (41.6 percent of 60 percent of available port bandwidth)

The 25 that is the numerator in $25/(15+20+25)$ is the percentage of bandwidth allocated to the fabric_fcset_be by the transmit rate weight. The (15+20+25) in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the three fabric fc-sets in the unicast class group.

The 60 percent represents 60 percent of the port bandwidth available after serving strict-high priority traffic. If no strict-high priority traffic is on the system, the equation results in the following bandwidth allocation to the fabric_fcset_be:

$25/(15+20+25)$ percent \times (0.6 \times 40 gigabits) = approximately 9.98 gigabits

- fabric_fcset_noloss1

15/(15+20+25) percent of 60 percent of the available port bandwidth (25 percent of 60 percent of available port bandwidth)

The 15 that is the numerator in **15**/(15+20+25) is the percentage of bandwidth allocated to the fabric_fcset_noloss1 by the transmit rate weight. The (15+20+25) in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the three fabric fc-sets in the unicast class group.

The 60 percent represents 60 percent of the port bandwidth available after serving strict-high priority traffic. If no strict-high priority traffic is on the system, the equation results in the following bandwidth allocation to the fabric_fcset_noloss1:

$$15/(15+20+25) \text{ percent} \times (0.6 \times 40 \text{ gigabits}) = \text{approximately } 6 \text{ gigabits}$$

- fabric_fcset_noloss2

20/(15+20+25) percent of 60 percent of the available port bandwidth (33.3 percent of 60 percent of available port bandwidth)

The 20 that is the numerator in **20**/(15+20+25) is the percentage of bandwidth allocated to the fabric_fcset_noloss2 by the transmit rate weight. The (15+20+25) in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the three fabric fc-sets in the unicast class group.

The 60 percent represents 60 percent of the port bandwidth available after serving strict-high priority traffic. If no strict-high priority traffic is on the system, the equation results in the following bandwidth allocation to the fabric_fcset_noloss2:

$$20/(15+20+25) \text{ percent} \times (0.6 \times 40 \text{ gigabits}) = \text{approximately } 7.99 \text{ gigabits}$$

The multidestination class group bandwidth is divided among its two fabric fc-sets:

- fabric_fcset_multicast1

10/(10+30) percent of 40 percent of the available port bandwidth (25 percent of 40 percent of available port bandwidth)

The 10 that is the numerator in **10**/(10+30) is the percentage of bandwidth allocated to the fabric_fcset_multicast1 by the transmit rate weight. The (10+30) in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the two fabric fc-sets in the multidestination class group.

The 40 percent represents 40 percent of the port bandwidth available after serving strict-high priority traffic. If no strict-high priority traffic is on the system, the equation results in the following bandwidth allocation to the fabric_fcset_multicast1:

$$10/(10+30) \text{ percent} \times (0.4 \times 40 \text{ gigabits}) = \text{approximately } 4 \text{ gigabits}$$

- fabric_fcset_multicast2

30/(10+30) percent of 40 percent of the available port bandwidth (75 percent of 40 percent of available port bandwidth)

The 30 that is the numerator in **30**/(10+30) is the percentage of bandwidth allocated to the fabric_fcset_multicast2 by the transmit rate weight. The (10+30) in the denominator sums the percentage of bandwidth (transmit rate weights) allocated to each of the two fabric fc-sets in the multideestination class group.

The 40 percent represents 40 percent of the port bandwidth available after serving strict-high priority traffic. If no strict -high priority traffic is on the system, the equation results in the following bandwidth allocation to the fabric_fcset_multicast2:

30/(10+30) percent x (0.4 x 40 gigabits) = approximately 12 gigabits

Configuring Scheduling on Node Device Fabric Interfaces

If you do not want to use default CoS scheduling on Node device fabric interfaces, you can configure two-tier hierarchical scheduling (ETS) the same way that you configure ETS on Node device access interfaces.

Similarities Between Node Device Fabric Interface and Access Interface Scheduling

Configuring scheduling on a Node device fabric interface is similar to configuring scheduling on an access interface in many ways. In both cases, you configure:

- Schedulers to specify the output scheduling for forwarding class traffic
- Scheduler maps to map schedulers to forwarding classes
- Forwarding classes (or use the default forwarding classes)
- Forwarding class sets (groups of forwarding classes that require similar CoS treatment)
- A separate fc-set for strict-high priority traffic (an fc-set cannot contain a mix of strict-high priority traffic and traffic with a different priority)
- Traffic control profiles to specify the output scheduling for fc-sets
- Traffic control profile and fc-set mapping to interfaces

On Node device fabric interfaces, you configure ETS in the same way and ETS works the same way as on Node device access interfaces

In addition, strict-high priority queues are served first, and then the remaining port bandwidth is allocated to other traffic.

Differences Between Node Device Fabric Interface and Access Interface Scheduling

Configuring scheduling on a Node device fabric interface differs from configuring scheduling on an access interface in several ways. On fabric interfaces:

- You cannot configure classifiers.
- You cannot configure congestion notification profiles (flow control is applied automatically to lossless forwarding classes).
- You specify the interface name differently.

Congestion Management

The Interconnect device is a shared component for all of the connected Node devices. Configuring scheduling on the external fabric interfaces (fte) and the internal Clos fabric interfaces (bfte) enables you to ensure predictable bandwidth usage for traffic flows across the Interconnect device.

Although minimal congestion is expected on the 40-Gbps fabric interfaces, you should configure congestion management to control packet drop during periods of congestion.

Lossy (Best Effort) Unicast Traffic

For unicast traffic that does not require lossless treatment, configure drop profiles (the standard Junos OS packet drop mechanism) to control packet drop during periods of congestion. (Drop profiles are not applied to multidestination traffic.)

A drop profile sets weighted random early detection (WRED) thresholds for dropping packets under different levels of congestion. Congestion levels for packet drop thresholds are fill levels of the output queue. When the output queue fills to a configured threshold, packet drop begins at the configured drop rate. When the output queue fills to a second configured threshold, packet drop reaches the configured maximum drop rate. You can apply different drop profiles to different types of traffic to achieve the desired pattern of packet loss during periods of congestion.

We recommend that you configure a relatively aggressive drop profile for traffic with a high loss priority and a less aggressive drop profile for traffic with a lower loss priority.

To create a drop profile and apply it to traffic of a certain loss priority:

1. Set loss priorities (low, medium-high, high) for different types of lossy unicast traffic when you configure classifiers and apply them to Node device access interfaces.
2. For each loss priority, configure at least one drop profile to define the WRED packet drop probability at different queue fill level thresholds. Create a more aggressive drop profile for traffic with a high loss priority, and progressively less aggressive drop profiles for traffic with medium-high and low loss priorities.

3. As part of scheduler configuration, configure a drop profile map, which maps a drop profile to a loss priority. A scheduler drop profile map can include mapping each loss priority to a drop profile, so you can specify different drop profiles for different traffic loss priorities in one scheduler. The scheduler uses the configured drop profile map to apply different drop profiles to traffic of different loss priorities, and thus control packet drop during periods of congestion.

Lossless Traffic

Do not configure drop profiles for lossless traffic. If you intend to map a scheduler to a lossless fabric fc-set, do not configure a drop profile for that scheduler.

The QFabric system automatically applies flow control that is similar to priority-based flow control (PFC) to traffic in lossless fabric fc-sets to prevent packet loss. In addition, you must enable PFC on the IEEE 802.1p code points for the lossless traffic at the Node device ingress interface to support lossless transport across the QFabric system. (You should also configure PFC across the Ethernet network to support lossless transport across the rest of the network.)

Multidestination Traffic

Drop profiles are not supported for multidestination traffic. Do not configure a drop profile in schedulers that you want to use for multidestination traffic.

RELATED DOCUMENTATION

[Understanding CoS Fabric Forwarding Class Sets | 396](#)

[Understanding CoS Output Queue Schedulers | 186](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 223](#)

[Understanding Default CoS Scheduling and Classification | 95](#)

[Understanding Default CoS Settings | 33](#)

[Example: Configuring CoS Scheduling Across the QFabric System | 353](#)

[Example: Configuring Queue Schedulers | 198](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 230](#)

[Example: Configuring WRED Drop Profiles | 300](#)

[Example: Configuring Drop Profile Maps | 307](#)

Example: Configuring CoS Scheduling Across the QFabric System

IN THIS SECTION

- [Requirements | 353](#)
- [Overview | 353](#)
- [Configuration | 365](#)
- [Verification | 381](#)

If you do not want to use the default class of service (CoS) scheduling of traffic across the QFabric system, then in addition to configuring CoS on Node device access interfaces, you can configure two-tier hierarchical scheduling on the fabric interfaces of a QFabric system. Configuring CoS on the fabric interfaces provides more control over class of service (CoS) across the QFabric system and helps to ensure predictable bandwidth consumption across the fabric path.

This topic describes:

Requirements

This example uses the following hardware and software components:

- Juniper Networks QFabric System with two Juniper Networks QFX3500 Node devices
- Junos OS Release 12.3 or later for the QFX Series

Overview

IN THIS SECTION

- [Topology | 356](#)

Configuring CoS across the QFabric system enables you to control scheduling resources as traffic passes through each type of interface. You can configure CoS on the following QFabric system interface types:

- Node device access interfaces (xe interfaces)—Schedule traffic on the output queues of the 10-Gigabit Ethernet access ports, using standard Node device CoS scheduling configuration

components, as described elsewhere in the QFX Series documentation. You can configure different scheduling for different ports and queues.

- Node device fabric interfaces (fte interfaces)—Schedule traffic on the output queues of the 40-Gbps fabric interfaces that connect a Node device to a QFX3008-I or a QFX3600-I Interconnect device using standard Node device CoS scheduling configuration components. You can configure different scheduling for different interfaces and output queues.
- Interconnect device fabric interfaces (fte interfaces)—Schedule traffic on the output queues of the 40-Gbps fabric interfaces that connect an Interconnect device to a Node device. You can configure different scheduling for different interfaces and fabric forwarding class sets (fabric fc-sets).
- Interconnect device internal Clos fabric interfaces (bfte interfaces)—Schedule traffic on the internal 40-Gbps Clos fabric interfaces that connect the three stages of the Clos fabric within the Interconnect device. You can configure one Clos fabric interface scheduler, which is applied to all of the internal Clos fabric interfaces. You cannot configure different schedulers for different Clos fabric interfaces.

This example shows you how to configure hierarchical port scheduling across the QFabric, including the configuration of Node device access interfaces, Node device fabric interfaces, Interconnect device fabric interfaces, and internal Interconnect device Clos fabric interfaces.

Configuring CoS on Interconnect device fabric interfaces differs from configuring CoS on Node device interfaces because the Interconnect device is a shared infrastructure that supports traffic from multiple Node devices and multiple Node device CoS configurations. Take the amounts and types of traffic traversing the Interconnect device into account when you configure CoS on Interconnect device interfaces.

Configuring scheduling across the QFabric system entails configuring interfaces on Node devices and Interconnect devices. You configure some or all of the following CoS components on each interface, depending upon the interface type (access, Node fabric, Interconnect fabric, or Interconnect Clos fabric):

- Mapping forwarding classes to priorities (IEEE 802.1p code points) and queues, and configuring lossless forwarding classes
- Defining fc-sets (priority groups)
- Defining drop profiles
- Defining schedulers
- Mapping forwarding classes to schedulers (scheduler map on Node devices, fabric scheduler map on Interconnect devices)
- Defining traffic control profiles

- Configuring a congestion notification profile to enable priority-based flow control (PFC) on lossless forwarding classes (priorities) (Node device access interfaces only)
- Applying congestion notification profiles to interfaces (Node device access interfaces only)
- Assigning fc-sets and traffic control profiles to interfaces (Node device interfaces only) or assigning fabric scheduler maps to interfaces (Interconnect device interfaces only)

NOTE: This example uses the default behavior aggregate classifiers on the Node device access interfaces. Classifiers are not applied to fabric interfaces. Although packet classification is not scheduling, it controls the forwarding class mapping to IEEE 802.1p priorities, and the loss priorities to which packets are mapped when they enter Node device access ports.

When you plan port bandwidth scheduling for priority groups (fc-sets on Node devices and class groups on Interconnect devices) and priorities (forwarding classes on Node devices and fabric fc-sets on Interconnect devices), take into account:

- The amounts and types of traffic you expect to traverse the Node device interfaces
- The amounts and types of aggregated traffic from all of the connected Node devices that you expect to traverse the Interconnect device interfaces
- The mapping of priorities into priority groups. Traffic that requires similar treatment usually belongs in the same priority group. To do this on Node devices, place forwarding classes that require similar bandwidth, loss priority, and other characteristics in the same fc-set. For example, you can map all types of best-effort traffic forwarding classes into one fc-set. On Interconnect devices, the default mapping of fabric fc-sets to class groups defines priority group membership and is not user-configurable.
- How much of the port bandwidth you want to allocate to each priority group and to each of the priorities in each priority group. The following considerations apply to bandwidth allocation:
 - Estimate how much traffic you expect in each priority's output queue (forwarding class on Node devices and fabric fc-set on Interconnect devices) and how much traffic you expect in each priority group (fc-set on Node devices and class group on Interconnect devices). The priority group traffic is the aggregated amount of traffic in the priorities that belong to the priority group.
 - On Node devices, the combined minimum guaranteed bandwidth of the priorities in a priority group should not exceed the minimum guaranteed bandwidth (guaranteed rate) of the priority group. (On Interconnect devices, class group bandwidth is derived from the bandwidth of the member fabric fc-sets, so the sum of the priority bandwidths cannot exceed the priority group bandwidth.) The transmit rate scheduler parameter defines the minimum guaranteed bandwidth for priorities (forwarding classes and fabric fc-sets). Scheduler maps associate schedulers with

forwarding classes (Node devices) and fabric scheduler maps associate schedulers with fabric fc-sets (Interconnect devices).

- The combined minimum guaranteed bandwidth of all of the priority groups on an interface should not exceed the interface's total bandwidth.

Topology

Figure 15 on page 356 shows the network topology used in this example.

Figure 15: Network Topology for Scheduling Across the QFabric System

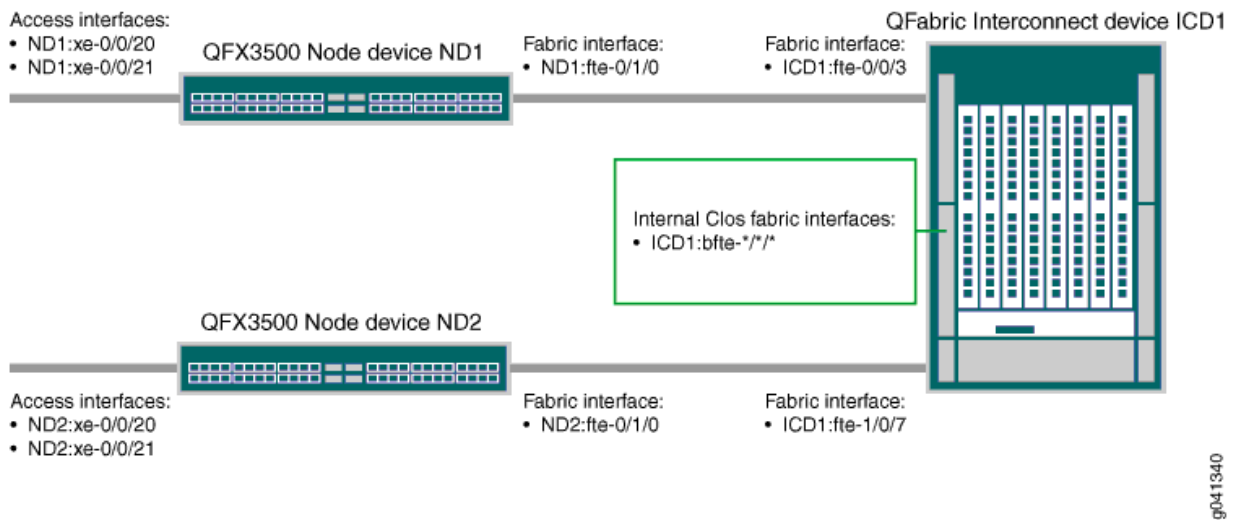


Table 76 on page 357 and Table 77 on page 362 describe the scheduling configuration components on the Node device and the Interconnect device.

To simplify Node device configuration, this example uses the same scheduling configuration on the access interfaces and the fabric interfaces of both QFabric Node devices. This is possible because the scheduler (forwarding-class scheduling) and traffic control profile (fc-set scheduling) rates are specified as percentages of bandwidth instead of as absolute values, so the schedulers and traffic control profiles utilize the port bandwidth in the same way regardless of the absolute amount of available bandwidth. If you want to treat traffic differently on different interfaces or on different interface types, you can configure different schedulers and traffic control profiles and apply them to the interfaces.

Table 76 on page 357 shows the scheduling configuration components for Node device interfaces.

Table 76: Components of the QFabric Node Device Hierarchical Port Scheduling Configuration Topology

Scheduling Component	Settings
Hardware	Two QFX3500 Node devices in a QFabric system
Forwarding classes	<p>This example uses five forwarding classes:</p> <ul style="list-style-type: none"> • best-effort • fcoe • no-loss • network-control • mcast <p>This example uses the default configuration for three forwarding classes (best-effort, network-control, and mcast). Best-effort traffic is classified into low loss priority and high loss priority by IEEE 802.1p classifiers at the Node device ingress interfaces.</p> <p>The two lossless forwarding classes (fcoe and no-loss) are configured as lossless forwarding classes:</p> <ul style="list-style-type: none"> • fcoe—Mapped to queue 3 with the no-loss parameter specified • no-loss—Mapped to queue 4 with the no-loss parameter specified <p>NOTE: Starting with Junos OS Release 12.3, you must include the no-loss parameter in the forwarding class configuration for forwarding classes that you want to be lossless. In Junos OS Release 12.3, all default forwarding classes, including the fcoe and no-loss forwarding classes, are lossy forwarding classes by default and must be explicitly configured as lossless to receive lossless CoS treatment. This is a change from lossless forwarding class configuration in earlier releases.</p>

Table 76: Components of the QFabric Node Device Hierarchical Port Scheduling Configuration Topology (Continued)

Scheduling Component	Settings
Forwarding class sets (priority groups)	<p>best-effort-pg—contains the forwarding classes best-effort and network-control</p> <p>no-loss-pg—contains the forwarding classes fcoe and no-loss</p> <p>multidestination-pg—contains the forwarding class mcast</p>
<p>Drop profiles</p> <p>NOTE: Lossless traffic (fcoe and no-loss forwarding classes) and multidestination traffic do not use drop profiles</p>	<p>This example uses the following drop profiles for lossy traffic classes:</p> <ul style="list-style-type: none"> Best-effort unicast traffic with low packet loss priority: <ul style="list-style-type: none"> Name—dp-be-low Drop start point—25% Drop end point—50% Maximum drop rate—80% Best-effort traffic unicast with high packet loss priority: <ul style="list-style-type: none"> Name—dp-be-high Drop start point—10% Drop end point—40% Maximum drop rate—100% Network-control traffic: <ul style="list-style-type: none"> Name—dp-nc Drop start point—75% Drop end point—100% Maximum drop rate—50%

Table 76: Components of the QFabric Node Device Hierarchical Port Scheduling Configuration Topology (Continued)

Scheduling Component	Settings
Queue (forwarding class) schedulers	<p>Schedulers configure the bandwidth characteristics of forwarding classes, which are mapped to output queues and to IEEE 802.1p CoS priorities.</p> <ul style="list-style-type: none"> Best-effort traffic scheduler: <ul style="list-style-type: none"> Name—be-sched Transmit rate (minimum guaranteed bandwidth)—90% Shaping rate (maximum bandwidth)—100% Priority—low Drop profiles—dp-be-low and dp-be-high Network-control traffic scheduler: <ul style="list-style-type: none"> Name—nc-sched Transmit rate—10% Shaping rate—100% Priority—low Drop profile—dp-nc FCoE traffic scheduler: <ul style="list-style-type: none"> Name—fcoe-sched Transmit rate—60% Shaping rate—100% Priority—low Drop profile—None No-loss traffic scheduler: <ul style="list-style-type: none"> Name—nl-sched Transmit rate—40% Shaping rate—100% Priority—low Drop profile—None Multidestination traffic scheduler: <ul style="list-style-type: none"> Name—mcast-sched Transmit rate—100% Shaping rate—100% Priority—low Drop profile—None <p>NOTE: If you want to specify absolute values instead of percentages for the transmit rate and the shaping rate, you should create</p>

Table 76: Components of the QFabric Node Device Hierarchical Port Scheduling Configuration Topology (Continued)

Scheduling Component	Settings
	<p>separate schedulers for access and fabric interfaces, because access interfaces are 10-Gigabit Ethernet interfaces and fabric interfaces are 40-Gbps interfaces.</p>
Forwarding class to scheduler mapping	<ul style="list-style-type: none"> Best-effort traffic scheduler map: Name—be-map Mapping—forwarding class best-effort to scheduler be-sched, forwarding class network-control to scheduler nc-sched Lossless traffic scheduler map: Name—nl-map Mapping—forwarding class fcoe to scheduler fcoe-sched, forwarding class no-loss to scheduler nl-sched Multidestination traffic scheduler map: Name—mcast-map Mapping—forwarding class mcast to scheduler mcast-sched
Priority group (fc-set) traffic control profiles	<p>Traffic control profiles configure the bandwidth for fc-sets (priority groups) and control the amount of port bandwidth allocated to the forwarding classes in the fc-sets.</p> <ul style="list-style-type: none"> Best-effort traffic control profile: Name—be-tcp Guaranteed rate (minimum guaranteed bandwidth)—25% Shaping rate (maximum bandwidth)—100% Scheduler map—be-map Lossless traffic control profile: Name—nl-tcp Guaranteed rate—50% Shaping rate—100% Scheduler map—nl-map Multidestination traffic control profile: Name—mcast-tcp Guaranteed rate—25% Shaping rate—100% Scheduler map—mcast-map

Table 76: Components of the QFabric Node Device Hierarchical Port Scheduling Configuration Topology (Continued)

Scheduling Component	Settings
Hierarchical scheduling (fc-sets and traffic control profiles) association with interfaces	<p>Apply the fc-sets and traffic control profiles to the interfaces of both Node devices:</p> <ul style="list-style-type: none"> Access interfaces—ND1:xe-0/0/20, ND1:xe-0/0/21, ND2:xe-0/0/20, ND2:xe-0/0/21 Fabric interfaces—ND1:fte-0/1/0, ND2:fte-0/1/0
PFC (access interfaces only; do not apply PFC to fabric interfaces)	<p>Code points:</p> <p>011—fcoe forwarding class traffic priority</p> <p>010—no-loss forwarding class traffic priority</p> <p>Congestion notification profile name—nl-cnp</p> <p>Enabled on interfaces: ND1:xe-0/0/20, ND1:xe-0/0/21, ND2:xe-0/0/20, and ND2:xe-0/0/21</p>

To simplify Interconnect device configuration, this example uses the same scheduling configuration on the fabric interfaces and the Clos fabric interfaces. If you want to treat traffic differently on different fabric interfaces or on different fabric interface types, you can configure different fabric schedulers, map them to fabric fc-sets, and apply them to the interfaces. (You can apply different mappings of schedulers to fabric fc-sets on different interfaces.)

NOTE: On Interconnect devices, the network-control forwarding class is mapped by default to the strict-high priority fabric fc-set (fabric_fcset_strict_high). The strict-high priority fabric fc-set receives all of the port bandwidth it needs to service strict-high priority traffic. You can configure a scheduler with a shaping rate (maximum bandwidth) and a drop profile to limit the bandwidth available to the strict-high priority fabric fc-set, if desired. The available fabric port bandwidth for all other traffic in all other fabric fc-sets is the bandwidth that remains after the interface services the strict-high priority traffic.

Table 77 on page 362 shows the scheduling configuration components for Interconnect device interfaces:

Table 77: Components of the QFabric Interconnect Device Hierarchical Port Scheduling Configuration Topology

Fabric Scheduling Component	Settings
Hardware	One QFabric Interconnect device connected to two QFX3500 Node devices in a QFabric system
Forwarding classes	<p>Interconnect devices use the forwarding classes defined on the connected Node devices. The forwarding classes are mapped by default to fabric fc-sets on the Interconnect device.</p> <p>NOTE: If you do not want to use the default forwarding class to fabric fc-set mapping, you can configure the mapping. Forwarding class to fabric fc-set mapping is global and applies to all traffic that crosses the Interconnect device.</p>
Fabric fc-sets	<p>This example uses four of the default fabric fc-sets, with the default mapping of forwarding classes to fabric fc-sets:</p> <ul style="list-style-type: none"> • fabric_fcset_be (includes the best-effort forwarding class) • fabric_fcset_noloss1 (includes the fcoe forwarding class) • fabric_fcset_noloss2 (includes the no-loss forwarding class) • fabric_fcset_multicast1 (includes the mcast forwarding class)
Class groups (priority groups)	The three default class groups and fabric fc-set membership in the class groups are not user-configurable.

Table 77: Components of the QFabric Interconnect Device Hierarchical Port Scheduling Configuration Topology (Continued)

Fabric Scheduling Component	Settings
<p>Drop profiles</p> <p>NOTE: Lossless traffic (fabric_fcset_noloss1 and fabric_fcset_noloss2) multideestination traffic do not use drop profiles</p>	<p>This example uses the following drop profiles for lossy traffic classes:</p> <ul style="list-style-type: none"> • Best-effort unicast traffic with low packet loss priority: <ul style="list-style-type: none"> Name—fab-dp-be-low Drop start point—20% Drop end point—50% Maximum drop rate—80% • Best-effort unicast traffic with high packet loss priority: <ul style="list-style-type: none"> Name—fab-dp-be-high Drop start point—5% Drop end point—35% Maximum drop rate—100%

Table 77: Components of the QFabric Interconnect Device Hierarchical Port Scheduling Configuration Topology (Continued)

Fabric Scheduling Component	Settings
Queue (fabric fc-set) fabric schedulers	<p>Schedulers configure the bandwidth for fabric fc-sets, which are mapped to output queues and to IEEE 802.1p CoS priorities.</p> <p>The sum of the minimum guaranteed bandwidths (transmit rates) of each fabric fc-set in a class group equals the total minimum guaranteed port bandwidth of the class group. The sum of all of the fabric fc-set transmit rates in all of the class groups equals the percentage of available port bandwidth allocated to the class groups. The sum of all of the fabric fc-set transmit rates must be less than or equal to 100 percent.</p> <ul style="list-style-type: none"> • Best-effort traffic scheduler: <ul style="list-style-type: none"> Name—fab-be-sched Transmit rate—25% Shaping rate—100% Drop profiles—fab-dp-be-low, fab-dp-be-high • FCoE traffic scheduler: <ul style="list-style-type: none"> Name—fab-fcoe-sched Transmit rate—30% Shaping rate—100% Drop profile—None • No-loss traffic scheduler: <ul style="list-style-type: none"> Name—fab-nl-sched Transmit rate—25% Shaping rate—100% Drop profile—None • Multidestination traffic scheduler: <ul style="list-style-type: none"> Name—fab-mcast-sched Transmit rate—20% Shaping rate—100% Drop profile—None

Table 77: Components of the QFabric Interconnect Device Hierarchical Port Scheduling Configuration Topology (Continued)

Fabric Scheduling Component	Settings
Fabric fc-set to fabric forwarding class set scheduler mapping	<ul style="list-style-type: none"> Best-effort traffic fabric scheduler mapping: Name—fab-traffic-map Mapping—fabric_fcset_be to scheduler fab-be-sched FCoE traffic fabric scheduler mapping: Name—fab-traffic-map Mapping—fabric_fcset_noloss1 to scheduler fab-fcoe-sched No-loss traffic fabric scheduler mapping: Name—fab-traffic-map Mapping—fabric_fcset_noloss2 to scheduler fab-nl-sched Multidestination traffic fabric scheduler mapping: Name—fab-traffic-map Mapping—fabric_fcset_mcast1 to scheduler fab-mcast-sched
Applying hierarchical scheduling (fabric scheduler map) to interfaces	<p>Fabric interfaces: ICD1:fte-0/0/3, ICD1:fte-1/0/7</p> <p>Clos fabric interfaces: ICD1:bftc-*/*/*</p>

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 366](#)
- [Configuring QFX3500 Node Devices ND1 and ND2 | 369](#)
- [Configuring QFX3500 Interconnect Device ICD1 | 374](#)
- [Results | 376](#)

The configuration example is split into two parts, one part for Node device scheduling configuration and one part for Interconnect device scheduling configuration. Although this example uses the same scheduling on Node device access and fabric interfaces, you can configure different schedulers for different interfaces. This example also uses the same scheduling on Interconnect device fabric and Clos fabric interfaces, and you can configure different schedulers for different interfaces.

To configure scheduling across a QFabric system, perform these tasks:

CLI Quick Configuration

Node device configuration: to quickly configure scheduling across a QFabric system, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI for Node devices ND1 and ND2 at the [edit] hierarchy level. In this example, we use identical scheduling and interfaces on Node devices ND1 and ND2 to simplify the configuration.

```
[edit class-of-service]
set forwarding-classes class fcoe queue-num 3 no-loss
set forwarding-classes class no-loss queue-num 4 no-loss
set forwarding-class-sets best-effort-pg class best-effort
set forwarding-class-sets best-effort-pg class network-control
set forwarding-class-sets no-loss-pg class fcoe
set forwarding-class-sets no-loss-pg class no-loss
set forwarding-class-sets multidestination-pg class mcast
set drop-profiles dp-be-low interpolate fill-level 25 fill-level 50 drop-probability 0 drop-
probability 80
set drop-profiles dp-be-high interpolate fill-level 10 fill-level 40 drop-probability 0 drop-
probability 100
set drop-profiles dp-nc interpolate fill-level 75 fill-level 100 drop-probability 0 drop-
probability 50
set schedulers be-sched priority low transmit-rate percent 90
set schedulers be-sched shaping-rate percent 100
set schedulers be-sched drop-profile-map loss-priority low protocol any drop-profile dp-be-low
set schedulers be-sched drop-profile-map loss-priority high protocol any drop-profile dp-be-high
set schedulers nc-sched priority low transmit-rate percent 10
set schedulers nc-sched shaping-rate percent 100
set schedulers nc-sched drop-profile-map loss-priority low protocol any drop-profile dp-nc
set schedulers fcoe-sched priority low transmit-rate percent 60
set schedulers fcoe-sched shaping-rate percent 100
set schedulers nl-sched priority low transmit-rate percent 40
set schedulers nl-sched shaping-rate percent 100
set schedulers mcast-sched priority low transmit-rate percent 100
set schedulers mcast-sched shaping-rate percent 100
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set scheduler-maps be-map forwarding-class network-control scheduler nc-sched
set scheduler-maps nl-map forwarding-class fcoe scheduler fcoe-sched
set scheduler-maps nl-map forwarding-class no-loss scheduler nl-sched
set scheduler-maps mcast-map forwarding-class mcast scheduler mcast-sched
```

```

set traffic-control-profiles be-tcp scheduler-map be-map guaranteed-rate percent 25
set traffic-control-profiles be-tcp shaping-rate percent 100
set traffic-control-profiles nl-tcp scheduler-map nl-map guaranteed-rate percent 50
set traffic-control-profiles nl-tcp shaping-rate percent 100
set traffic-control-profiles mcast-tcp scheduler-map mcast-map guaranteed-rate percent 25
set traffic-control-profiles mcast-tcp shaping-rate percent 100
set interfaces ND1:xe-0/0/20 forwarding-class-set best-effort-pg output-traffic-control-profile
be-tcp
set interfaces ND1:xe-0/0/20 forwarding-class-set noloop-pg output-traffic-control-profile nl-tcp
set interfaces ND1:xe-0/0/20 forwarding-class-set multidestination-pg output-traffic-control-
profile mcast-tcp
set interfaces ND1:xe-0/0/21 forwarding-class-set best-effort-pg output-traffic-control-profile
be-tcp
set interfaces ND1:xe-0/0/21 forwarding-class-set noloop-pg output-traffic-control-profile nl-tcp
set interfaces ND1:xe-0/0/21 forwarding-class-set multidestination-pg output-traffic-control-
profile mcast-tcp
set interfaces ND1:fte-0/1/0 forwarding-class-set best-effort-pg output-traffic-control-profile
be-tcp
set interfaces ND1:fte-0/1/0 forwarding-class-set noloop-pg output-traffic-control-profile nl-tcp
set interfaces ND1:fte-0/1/0 forwarding-class-set multidestination-pg output-traffic-control-
profile mcast-tcp
set interfaces ND2:xe-0/0/20 forwarding-class-set best-effort-pg output-traffic-control-profile
be-tcp
set interfaces ND2:xe-0/0/20 forwarding-class-set noloop-pg output-traffic-control-profile nl-tcp
set interfaces ND2:xe-0/0/20 forwarding-class-set multidestination-pg output-traffic-control-
profile mcast-tcp
set interfaces ND2:xe-0/0/21 forwarding-class-set best-effort-pg output-traffic-control-profile
be-tcp
set interfaces ND2:xe-0/0/21 forwarding-class-set noloop-pg output-traffic-control-profile nl-tcp
set interfaces ND2:xe-0/0/21 forwarding-class-set multidestination-pg output-traffic-control-
profile mcast-tcp
set interfaces ND2:fte-0/1/0 forwarding-class-set best-effort-pg output-traffic-control-profile
be-tcp
set interfaces ND2:fte-0/1/0 forwarding-class-set noloop-pg output-traffic-control-profile nl-tcp
set interfaces ND2:fte-0/1/0 forwarding-class-set multidestination-pg output-traffic-control-
profile mcast-tcp
set congestion-notification-profile nl-cnp input ieee-802.1 code-point 011 pfc
set congestion-notification-profile nl-cnp input ieee-802.1 code-point 100 pfc
set interfaces ND1:xe-0/0/20 congestion-notification-profile nl-cnp
set interfaces ND1:xe-0/0/21 congestion-notification-profile nl-cnp
set interfaces ND2:xe-0/0/20 congestion-notification-profile nl-cnp
set interfaces ND2:xe-0/0/21 congestion-notification-profile nl-cnp

```

Interconnect device configuration: to quickly configure scheduling across a QFabric system, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI for Interconnect device ICD1 at the [edit] hierarchy level. In this example, we use identical scheduling on the fabric interfaces and the Clos fabric interfaces to simplify the configuration.

NOTE: This configuration uses the default mapping of forwarding classes to fabric fc-sets.

```
[edit class-of-service]
set drop-profiles fab-dp-be-low interpolate fill-level 20 fill-level 50 drop-probability 0 drop-
probability 80
set drop-profiles fab-dp-be-high interpolate fill-level 5 fill-level 35 drop-probability 0 drop-
probability 100
set schedulers fab-be-sched transmit-rate percent 25
set schedulers fab-be-sched shaping-rate percent 100
set schedulers fab-be-sched drop-profile-map loss-priority low protocol any drop-profile fab-dp-
be-low
set schedulers fab-be-sched drop-profile-map loss-priority high protocol any drop-profile fab-dp-
be-high
set schedulers fab-fcoe-sched transmit-rate percent 30
set schedulers fab-fcoe-sched shaping-rate percent 100
set schedulers fab-nl-sched transmit-rate percent 25
set schedulers fab-nl-sched shaping-rate percent 100
set schedulers fab-mcast-sched transmit-rate percent 20
set schedulers fab-mcast-sched shaping-rate percent 100
set scheduler-map-forwarding-class-sets fab-traffic-map forwarding-class-set fabric_fcset_be
scheduler fab-be-sched
set scheduler-map-forwarding-class-sets fab-traffic-map forwarding-class-set
fabric_fcset_noloss1 scheduler fab-fcoe-sched
set scheduler-map-forwarding-class-sets fab-traffic-map forwarding-class-set
fabric_fcset_noloss2 scheduler fab-nl-sched
set scheduler-map-forwarding-class-sets fab-traffic-map forwarding-class-set fabric_fcset_mcast1
scheduler fab-mcast-sched
set interfaces ICD1:fte-0/0/3 scheduler-map-forwarding-class-sets fab-traffic-map
set interfaces ICD1:fte-1/0/7 scheduler-map-forwarding-class-sets fab-traffic-map
set interfaces ICD1:bft-*/*/ scheduler-map-forwarding-class-sets fab-traffic-map
```

Configuring QFX3500 Node Devices ND1 and ND2

Step-by-Step Procedure

To perform a step-by-step configuration of lossless forwarding classes, forwarding class sets, drop profiles for lossy traffic, queue schedulers, traffic control profiles, access and fabric interfaces, and PFC:

1. Configure the two lossless forwarding classes (priorities):

```
[edit class-of-service]
user@switch# set forwarding-classes class fcoe queue-num 3 no-loss
user@switch# set forwarding-classes class no-loss queue-num 4 no-loss
```

2. Configure fc-sets (priority groups) to group forwarding classes (priorities) that require similar CoS treatment:

```
[edit class-of-service]
user@switch# set forwarding-class-sets best-effort-pg class best-effort
user@switch# set forwarding-class-sets best-effort-pg class network-control
user@switch# set forwarding-class-sets no-loss-pg class fcoe
user@switch# set forwarding-class-sets no-loss-pg class no-loss
user@switch# set forwarding-class-sets multidestination-pg class mcast
```

3. Configure the drop profile for the best-effort low loss-priority queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-be-low interpolate fill-level 25 fill-level 50 drop-
probability 0 drop-probability 80
```

4. Configure the drop profile for the best-effort high loss-priority queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-be-high interpolate fill-level 10 fill-level 40 drop-
probability 0 drop-probability 100
```

5. Configure the drop profile for the network-control queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-nc interpolate fill-level 75 fill-level 100 drop-
probability 0 drop-probability 50
```

6. Configure the scheduler that defines the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profiles for the best-effort queue:

```
[edit class-of-service]
user@switch# set schedulers be-sched priority low transmit-rate percent 90
user@switch# set schedulers be-sched shaping-rate percent 100
user@switch# set schedulers be-sched drop-profile-map loss-priority low protocol any drop-
profile dp-be-low
user@switch# set schedulers be-sched drop-profile-map loss-priority high protocol any drop-
profile dp-be-high
```

7. Configure the scheduler that defines the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profile for the network-control queue:

```
[edit class-of-service]
user@switch# set schedulers nc-sched priority low transmit-rate percent 10
user@switch# set schedulers nc-sched shaping-rate percent 100
user@switch# set schedulers nc-sched drop-profile-map loss-priority low protocol any drop-
profile dp-nc
```

8. Configure the scheduler that defines the minimum guaranteed bandwidth, priority, and maximum bandwidth for the FCoE queue:

```
[edit class-of-service]
user@switch# set schedulers fcoe-sched priority low transmit-rate percent 60
user@switch# set schedulers fcoe-sched shaping-rate percent 100
```

9. Configure the scheduler that defines the minimum guaranteed bandwidth, priority, and maximum bandwidth for the no-loss queue:

```
[edit class-of-service]
user@switch# set schedulers nl-sched priority low transmit-rate percent 40
user@switch# set schedulers nl-sched shaping-rate percent 100
```

10. Configure the scheduler that defines the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profile for the mcast queue:

```
[edit class-of-service]
user@switch# set schedulers mcast-sched priority low transmit-rate percent 100
user@switch# set schedulers mcast-sched shaping-rate percent 100
```

11. Map the schedulers to the appropriate forwarding classes:

```
[edit class-of-service]
user@switch# set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
user@switch# set scheduler-maps be-map forwarding-class network-control scheduler nc-sched
user@switch# set scheduler-maps nl-map forwarding-class fcoe scheduler fcoe-sched
user@switch# set scheduler-maps nl-map forwarding-class no-loss scheduler nl-sched
user@switch# set scheduler-maps mcast-map forwarding-class mcast scheduler mcast-sched
```

12. Define the traffic control profile for the best-effort priority group (queue to scheduler mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
user@switch# set traffic-control-profiles be-tcp scheduler-map be-map guaranteed-rate
percent 25
user@switch# set traffic-control-profiles be-tcp shaping-rate percent 100
```

13. Define the traffic control profile for the guaranteed delivery priority group (queue to scheduler mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
user@switch# set traffic-control-profiles nl-tcp scheduler-map nl-map guaranteed-rate
```

```
percent 50
```

```
user@switch# set traffic-control-profiles nl-tcp shaping-rate percent 100
```

14. Define the traffic control profile for the multideestination priority group (queue to scheduler mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
```

```
user@switch# set traffic-control-profiles mcast-tcp scheduler-map mcast-map guaranteed-rate  
percent 25
```

```
user@switch# set traffic-control-profiles mcast-tcp shaping-rate percent 100
```

15. Apply the three forwarding class sets and the appropriate traffic control profiles to the Node device ND1 access interfaces and fabric interface:

```
[edit class-of-service]
```

```
user@switch# set interfaces ND1:xe-0/0/20 forwarding-class-set best-effort-pg output-  
traffic-control-profile be-tcp
```

```
user@switch# set interfaces ND1:xe-0/0/20 forwarding-class-set noloss-pg output-traffic-  
control-profile nl-tcp
```

```
user@switch# set interfaces ND1:xe-0/0/20 forwarding-class-set multideestination-pg output-  
traffic-control-profile mcast-tcp
```

```
user@switch# set interfaces ND1:xe-0/0/21 forwarding-class-set best-effort-pg output-  
traffic-control-profile be-tcp
```

```
user@switch# set interfaces ND1:xe-0/0/21 forwarding-class-set noloss-pg output-traffic-  
control-profile nl-tcp
```

```
user@switch# set interfaces ND1:xe-0/0/21 forwarding-class-set multideestination-pg output-  
traffic-control-profile mcast-tcp
```

```
user@switch# set interfaces ND1:fte-0/1/0 forwarding-class-set best-effort-pg output-  
traffic-control-profile be-tcp
```

```
user@switch# set interfaces ND1:fte-0/1/0 forwarding-class-set noloss-pg output-traffic-  
control-profile nl-tcp
```

```
user@switch# set interfaces ND1:fte-0/1/0 forwarding-class-set multideestination-pg output-  
traffic-control-profile mcast-tcp
```

16. Apply the three forwarding class sets and the appropriate traffic control profiles to the Node device ND2 access interfaces and fabric interface:

```
[edit class-of-service]
```

```
user@switch# set interfaces ND2:xe-0/0/20 forwarding-class-set best-effort-pg output-
```

```

traffic-control-profile be-tcp
user@switch# set interfaces ND2:xe-0/0/20 forwarding-class-set no-loss-pg output-traffic-
control-profile nl-tcp
user@switch# set interfaces ND2:xe-0/0/20 forwarding-class-set multideestination-pg output-
traffic-control-profile mcast-tcp
user@switch# set interfaces ND2:xe-0/0/21 forwarding-class-set best-effort-pg output-
traffic-control-profile be-tcp
user@switch# set interfaces ND2:xe-0/0/21 forwarding-class-set no-loss-pg output-traffic-
control-profile nl-tcp
user@switch# set interfaces ND2:xe-0/0/21 forwarding-class-set multideestination-pg output-
traffic-control-profile mcast-tcp
user@switch# set interfaces ND2:fte-0/1/0 forwarding-class-set best-effort-pg output-
traffic-control-profile be-tcp
user@switch# set interfaces ND2:fte-0/1/0 forwarding-class-set no-loss-pg output-traffic-
control-profile nl-tcp
user@switch# set interfaces ND2:fte-0/1/0 forwarding-class-set multideestination-pg output-
traffic-control-profile mcast-tcp

```

17. Configure a congestion notification profile to enable PFC on the FCoE and no-loss queue IEEE 802.1 code points:

```

[edit class-of-service]
user@switch# set congestion-notification-profile nl-cnp input ieee-802.1 code-point 011 pfc
user@switch# set congestion-notification-profile nl-cnp input ieee-802.1 code-point 100 pfc

```

18. Apply the PFC configuration to the access interfaces on Node device ND1:

```

[edit class-of-service]
user@switch# set interfaces ND1:xe-0/0/20 congestion-notification-profile nl-cnp
set interfaces ND1:xe-0/0/21 congestion-notification-profile nl-cnp

```

19. Apply the PFC configuration to the access interfaces on Node device ND2:

```

[edit class-of-service]
user@switch# set interfaces ND2:xe-0/0/20 congestion-notification-profile nl-cnp
set interfaces ND2:xe-0/0/21 congestion-notification-profile nl-cnp

```


Configuring QFX3500 Interconnect Device ICD1

Step-by-Step Procedure

To perform a step-by-step configuration of drop profiles for lossy traffic, queue schedulers, and fabric and Clos fabric interfaces:

1. Configure the drop profile for the best-effort low loss-priority queue:

```
[edit class-of-service]
user@switch# set drop-profiles fab-dp-be-low interpolate fill-level 20 fill-level 50 drop-
probability 0 drop-probability 80
```

2. Configure the drop profile for the best-effort high loss-priority queue:

```
[edit class-of-service]
user@switch# set drop-profiles fab-dp-be-high interpolate fill-level 5 fill-level 35 drop-
probability 0 drop-probability 100
```

3. Configure the fabric scheduler that defines the minimum guaranteed bandwidth, maximum bandwidth, and drop profiles for the best-effort (fabric_fcset_be) queue:

```
[edit class-of-service]
user@switch# set schedulers fab-be-sched transmit-rate percent 25
user@switch# set schedulers fab-be-sched shaping-rate percent 100
user@switch# set schedulers fab-be-sched drop-profile-map loss-priority low protocol any drop-
profile fab-dp-be-low
user@switch# set schedulers fab-be-sched drop-profile-map loss-priority high protocol any
drop-profile fab-dp-be-high
```

4. Configure the fabric scheduler that defines the minimum guaranteed bandwidth and maximum bandwidth for the FCoE (fabric_fcset_noloss1) queue:

```
[edit class-of-service]
user@switch# set schedulers fab-fcoe-sched transmit-rate percent 30
user@switch# set schedulers fab-fcoe-sched shaping-rate percent 100
```

5. Configure the fabric scheduler that defines the minimum guaranteed bandwidth and maximum bandwidth for the no-loss (fabric_fcset_noloss2) queue:

```
[edit class-of-service]
user@switch# set schedulers fab-nl-sched transmit-rate percent 25
user@switch# set schedulers fab-nl-sched shaping-rate percent 100
```

6. Configure the fabric scheduler that defines the minimum guaranteed bandwidth, maximum bandwidth, and drop profile for the multidestination traffic (fabric_fcset_mcast1) queue:

```
[edit class-of-service]
user@switch# set schedulers fab-mcast-sched transmit-rate percent 20
user@switch# set schedulers fab-mcast-sched shaping-rate percent 100
```

7. Map the fabric schedulers to the appropriate fabric fc-sets in the fabric forwarding class scheduler map:

```
[edit class-of-service]
user@switch# set scheduler-map-forwarding-class-sets fab-traffic-map forwarding-class-set
fabric_fcset_be scheduler fab-be-sched
user@switch# set scheduler-map-forwarding-class-sets fab-traffic-map forwarding-class-set
fabric_fcset_noloss1 scheduler fab-fcoe-sched
user@switch# set scheduler-map-forwarding-class-sets fab-traffic-map forwarding-class-set
fabric_fcset_noloss2 scheduler fab-nl-sched
user@switch# set scheduler-map-forwarding-class-sets fab-traffic-map forwarding-class-set
fabric_fcset_mcast1 scheduler fab-mcast-sched
```

8. To configure scheduling on the interfaces, apply the scheduler map to the Interconnect device fabric interfaces and Clos fabric interfaces:

```
[edit class-of-service]
user@switch# set interfaces ICD1:fte-0/0/3 scheduler-map-forwarding-class-sets fab-traffic-map
user@switch# set interfaces ICD1:fte-1/0/7 scheduler-map-forwarding-class-sets fab-traffic-map
user@switch# set interfaces ICD1:bft-*/*/ scheduler-map-forwarding-class-sets fab-traffic-
map
```

Results

Display the results of the CoS configuration on QFX3500 Node devices ND1 and ND2. The system shows only the explicitly configured parameters; it does not show default parameters such as the classifier configuration or the default forwarding classes. In this example, the three lossy forwarding classes (best-effort, network-control, and mcast) are not shown because the example uses the default configuration for these forwarding classes. The results on both Node devices are similar, except the interface names are different because the interface names include the Node device name. The results below are for Node device ND1:

```
user@switch> show configuration class-of-service
```

```
drop-profiles {
  dp-be-low {
    interpolate {
      fill-level [ 25 50 ];
      drop-probability [ 0 80 ];
    }
  }
  dp-be-high {
    interpolate {
      fill-level [ 10 40 ];
      drop-probability [ 0 100 ];
    }
  }
  dp-nc {
    interpolate {
      fill-level [ 75 100 ];
      drop-probability [ 0 50 ];
    }
  }
}
forwarding-classes {
  class fcoe queue-num 3 no-loss;
  class no-loss queue-num 4 no-loss;
}
traffic-control-profiles {
  be-tcp {
    scheduler-map be-map;
    shaping-rate percent 100;
    guaranteed-rate percent 25;
  }
  nl-tcp {
```

```

        scheduler-map nl-map;
        shaping-rate percent 100;
        guaranteed-rate percent 50;
    }
    mcast-tcp {
        scheduler-map mcast-map;
        shaping-rate percent 100;
        guaranteed-rate percent 25;
    }
}
forwarding-class-sets {
    best-effort-pg {
        class best-effort;
        class network-control;
    }
    no-loss-pg {
        class fcoe;
        class no-loss;
    }
    multideestination-pg {
        class mcast;
    }
}
congestion-notification-profile {
    nl-cnp {
        input {
            ieee-802.1 {
                code-point 011 {
                    pfc;
                }
                code-point 100 {
                    pfc;
                }
            }
        }
    }
}
interfaces {
    ND1:xe-0/0/20 {
        congestion-notification-profile nl-cnp;
        forwarding-class-set {
            best-effort-pg {
                output-traffic-control-profile be-tcp;
            }
        }
    }
}

```

```

    }
    noloss-pg {
        output-traffic-control-profile nl-tcp;
    }
    multideestination-pg {
        output-traffic-control-profile mcast-tcp;
    }
}
}
ND1:xe-0/0/21 {
    congestion-notification-profile nl-cnp;
    forwarding-class-set {
        best-effort-pg {
            output-traffic-control-profile be-tcp;
        }
        noloss-pg {
            output-traffic-control-profile nl-tcp;
        }
        multideestination-pg {
            output-traffic-control-profile mcast-tcp;
        }
    }
}
}
ND1:fte-0/1/0 {
    forwarding-class-set {
        best-effort-pg {
            output-traffic-control-profile be-tcp;
        }
        noloss-pg {
            output-traffic-control-profile nl-tcp;
        }
        multideestination-pg {
            output-traffic-control-profile mcast-tcp;
        }
    }
}
}
scheduler-maps {
    be-map {
        forwarding-class best-effort scheduler be-sched;
        forwarding-class network-control scheduler nc-sched;
    }
    nl-map {

```

```

        forwarding-class fcoe scheduler fcoe-sched;
        forwarding-class no-loss scheduler nl-sched;
    }
    mcast-map {
        forwarding-class mcast scheduler mcast-sched;
    }
}
schedulers {
    be-sched {
        transmit-rate percent 90;
        shaping-rate percent 100;
        priority low;
        drop-profile-map loss-priority low protocol any drop-profile dp-be-low;
        drop-profile-map loss-priority high protocol any drop-profile dp-be-high;
    }
    fcoe-sched {
        transmit-rate percent 60;
        shaping-rate percent 100;
        priority low;
    }
    mcast-sched {
        transmit-rate percent 100;
        shaping-rate percent 100;
        priority low;
    }
    nc-sched {
        transmit-rate percent 10;
        shaping-rate percent 100;
        priority low;
        drop-profile-map loss-priority low protocol any drop-profile dp-nc;
    }
    nl-sched {
        transmit-rate percent 40;
        shaping-rate percent 100;
        priority low;
    }
}

```

Display the results of the CoS configuration on QFX3500 Interconnect device ICD1. The system shows only the explicitly configured parameters; it does not show default parameters:

```

user@switch> show configuration class-of-service
  drop-profiles {
    fab-dp-be-low {
      interpolate {
        fill-level [ 20 50 ];
        drop-probability [ 0 80 ];
      }
    }
    fab-dp-be-high {
      interpolate {
        fill-level [ 5 35 ];
        drop-probability [ 0 100 ];
      }
    }
  }
  interfaces {
    ICD1:fte-0/0/3 {
      scheduler-map-forwarding-class-sets fab-traffic-map;
    }
    ICD1:fte-1/0/7 {
      scheduler-map-forwarding-class-sets fab-traffic-map;
    }
    ICD1:bft-*/*/ {
      scheduler-map-forwarding-class-sets fab-traffic-map;
    }
  }
  scheduler-maps {
    fab-traffic-map {
      forwarding-class-set fabric_fcset_be scheduler fab-be-sched;
      forwarding-class-set fabric_fcset_noloss1 scheduler fab-fcoe-sched;
      forwarding-class-set fabric_fcset_noloss2 scheduler fab-nl-sched;
      forwarding-class-set fabric_fcset_mcast1 scheduler fab-mcast-sched;
    }
  }
  schedulers {
    fab-be-sched {
      transmit-rate percent 25;
      shaping-rate percent 100;
      drop-profile-map loss-priority low protocol any drop-profile fab-dp-be-low;
    }
  }

```

```

        drop-profile-map loss-priority high protocol any drop-profile fab-dp-be-high;
    }
    fab-fcoe-sched {
        transmit-rate percent 30;
        shaping-rate percent 100;
    }
    fab-nl-sched {
        transmit-rate percent 25;
        shaping-rate percent 100;
    }
    fab-mcast-sched {
        transmit-rate percent 20;
        shaping-rate percent 100;
    }
}

```

Verification

IN THIS SECTION

- [Verifying Lossless Forwarding Class Configuration on the Node Devices | 382](#)
- [Verifying Forwarding Class Set Configuration on the Node Devices | 382](#)
- [Verifying Drop Profile Configuration on the Node Devices | 383](#)
- [Verifying Drop Profile Configuration on the Interconnect Device | 384](#)
- [Verifying Queue Scheduler Configuration and Mapping on the Node Devices | 385](#)
- [Verifying Fabric Queue Scheduler Configuration and Mapping on the Interconnect Device | 388](#)
- [Verifying Traffic Control Profile Configuration on the Node Devices | 391](#)
- [Verifying That PFC Is Enabled on Lossless Queues on the Node Devices | 392](#)
- [Verifying Access and Fabric Interface Scheduling Configuration on the Node Devices | 393](#)
- [Verifying Fabric Interface Scheduling Configuration on the Interconnect Device | 394](#)

To verify that the hierarchical scheduling components have been created and are operating properly, perform these tasks:

Verifying Lossless Forwarding Class Configuration on the Node Devices

Purpose

On Node devices ND1 and ND2, verify that the two lossless forwarding classes (fcoe and no-loss) have been configured. The system shows only the explicitly configured forwarding classes, so the default configuration of the best-effort, network-control, and mcast forwarding classes is not shown.

Action

List the forwarding classes using the operational mode command `show configuration class-of-service forwarding-classes`:

```
user@switch> show configuration class-of-service forwarding-classes
class fcoe queue-num 3 no-loss;
class no-loss queue-num 4 no-loss;
```

Meaning

The `show configuration class-of-service forwarding-classes` command lists each of the configured forwarding classes, the queue to which the forwarding class is mapped, and whether the forwarding class has been configured to be lossless with the `no-loss` option. The command output shows that:

- Forwarding class `fcoe` maps to queue 3 and is configured as a lossless queue with the `no-loss` option
- Forwarding class `no-loss` maps to queue 4 and is configured as a lossless queue with the `no-loss` option

Verifying Forwarding Class Set Configuration on the Node Devices

Purpose

Verify that the correct forwarding classes belong to the appropriate fc-set.

Action

List the fc-sets on Node devices ND1 and ND2 using the operational mode command `show class-of-service forwarding-class-set`:

```
user@switch> show class-of-service forwarding-class-set
Forwarding class set: best-effort-pg, Type: normal-type, Forwarding class set index: 19907
```

Forwarding class	Index
best-effort	0
network-control	5

Forwarding class set: no-loss-pg, Type: normal-type, Forwarding class set index: 43700

Forwarding class	Index
fcoe	2
no-loss	3

Forwarding class set: multideestination-pg, Type: normal-type, Forwarding class set index: 60758

Forwarding class	Index
mcast	4

Meaning

The `show class-of-service forwarding-class-set` command lists all of the configured fc-sets (priority groups), the forwarding classes (priorities) that belong to each fc-set, and the internal index number of each fc-set. The command output shows that:

- The fc-set `best-effort-pg` includes the forwarding classes `best-effort` and `network-control`.
- The fc-set `no-loss-pg` includes the forwarding classes `fcoe` and `no-loss`.
- The fc-set `multideestination-pg` includes the forwarding class `mcast`.

Verifying Drop Profile Configuration on the Node Devices

Purpose

On Node devices ND1 and ND2, verify that the drop profiles `dp-be-low`, `dp-be-high`, and `dp-nc` are configured with the correct fill levels and drop probabilities.

Action

On Node devices ND1 and ND2, list the drop profiles using the operational mode command `show configuration class-of-service drop-profiles`:

```
user@switch> show configuration class-of-service drop-profiles
dp-be-low {
  interpolate {
    fill-level [ 25 50 ];
    drop-probability [ 0 80 ];
```

```

    }
}
dp-be-high {
    interpolate {
        fill-level [ 10 40 ];
        drop-probability [ 0 100 ];
    }
}
dp-nc {
    interpolate {
        fill-level [ 75 100 ];
        drop-probability [ 0 50 ];
    }
}

```

Meaning

The `show configuration class-of-service drop-profiles` command lists the drop profiles and their properties. The command output shows that there are three drop profiles configured, `dp-be-low`, `dp-be-high`, and `dp-nc`. The output also shows that:

- For `dp-be-low`, the drop start point (the first fill level) is when the queue is 25 percent filled, the drop end point (the second fill level) occurs when the queue is 50 percent filled, and the drop probability at the drop end point is 80 percent.
- For `dp-be-high`, the drop start point (the first fill level) is when the queue is 10 percent filled, the drop end point (the second fill level) occurs when the queue is 40 percent filled, and the drop probability at the drop end point is 100 percent.
- For `dp-nc`, the drop start point (the first fill level) is when the queue is 75 percent filled, the drop end point (the second fill level) occurs when the queue is 100 percent filled, and the drop probability at the drop end point is 50 percent.

Verifying Drop Profile Configuration on the Interconnect Device

Purpose

On Interconnect device ICD1, verify that drop profiles `fab-dp-be-low` and `fab-dp-be-high` are configured with the correct fill levels and drop probabilities.

Action

List the drop profiles using the operational mode command `show configuration class-of-service drop-profiles`:

```
user@switch> show configuration class-of-service drop-profiles
fab-dp-be-low {
    interpolate {
        fill-level [ 20 50 ];
        drop-probability [ 0 80 ];
    }
}
fab-dp-be-high {
    interpolate {
        fill-level [ 5 35 ];
        drop-probability [ 0 100 ];
    }
}
```

Meaning

The `show configuration class-of-service drop-profiles` command lists the drop profiles and their properties. The command output shows that there are two drop profiles configured, `fab-dp-be-low` and `fab-dp-be-high`. The output also shows that:

- For `fab-dp-be-low`, the drop start point (the first fill level) is when the queue is 20 percent filled, the drop end point (the second fill level) occurs when the queue is 50 percent filled, and the drop probability at the drop end point is 80 percent.
- For `fab-dp-be-high`, the drop start point (the first fill level) is when the queue is 5 percent filled, the drop end point (the second fill level) occurs when the queue is 35 percent filled, and the drop probability at the drop end point is 100 percent.

Verifying Queue Scheduler Configuration and Mapping on the Node Devices

Purpose

Verify that the Node device ND1 and ND2 queue schedulers are configured with the correct bandwidth parameters and priorities, mapped to the correct forwarding classes and queues, and mapped to the correct drop profiles.

Action

List the scheduler maps using the operational mode command `show class-of-service scheduler-map`:

```
user@switch> show class-of-service scheduler-map
Scheduler map: be-map, Index: 64023

Scheduler: be-sched, Forwarding class: best-effort, Index: 13005
  Transmit rate: 90 percent, Rate Limit: none, Buffer size: remainder,
  Buffer Limit: none, Priority: low
  Excess Priority: unspecified
  Shaping rate: 100 percent,
  Drop profiles:
    Loss priority  Protocol  Index  Name
    Low           any      55387  dp-be-low
    Medium high   any      1      <default-drop-profile>
    High          any      4369   dp-be-high

Scheduler: nc-sched, Forwarding class: network-control, Index: 45740
  Transmit rate: 10 percent, Rate Limit: none, Buffer size: remainder,
  Buffer Limit: none, Priority: low
  Excess Priority: unspecified
  Shaping rate: 100 percent,
  Drop profiles:
    Loss priority  Protocol  Index  Name
    Low           any      44207  dp-nc
    Medium high   any      1      <default-drop-profile>
    High          any      1      <default-drop-profile>

Scheduler map: nl-map, Index: 61447

Scheduler: fcoe-sched, Forwarding class: fcoe, Index: 37289
  Transmit rate: 60 percent, Rate Limit: none, Buffer size: remainder,
  Buffer Limit: none, Priority: low
  Excess Priority: unspecified
  Shaping rate: 100 percent,
  Drop profiles:
    Loss priority  Protocol  Index  Name
    Low           any      44207  <default-drop-profile>
    Medium high   any      1      <default-drop-profile>
    High          any      1      <default-drop-profile>
```

```
Scheduler: nl-sched, Forwarding class: no-loss, Index: 29359
  Transmit rate: 40 percent, Rate Limit: none, Buffer size: remainder,
  Buffer Limit: none, Priority: low
  Excess Priority: unspecified
  Shaping rate: 100 percent,
  Drop profiles:
    Loss priority  Protocol  Index  Name
    Low           any       44207  <default-drop-profile>
    Medium high   any       1      <default-drop-profile>
    High          any       1      <default-drop-profile>
```

```
Scheduler map: mcast-map, Index: 63239
```

```
Scheduler: mcast-sched, Forwarding class: mcast, Index: 29359
  Transmit rate: 100 percent, Rate Limit: none, Buffer size: remainder,
  Buffer Limit: none, Priority: low
  Excess Priority: unspecified
  Shaping rate: 100 percent,
  Drop profiles:
    Loss priority  Protocol  Index  Name
    Low           any       1      <default-drop-profile>
    Medium high   any       1      <default-drop-profile>
    High          any       1      <default-drop-profile>
```

Meaning

The `show class-of-service scheduler-map` command lists the three configured scheduler maps. For each scheduler map, the command output includes:

- The name of the scheduler map (Scheduler map field)
- The name of the scheduler (Scheduler field)
- The forwarding classes mapped to the scheduler (Forwarding class field)
- The minimum guaranteed queue bandwidth (Transmit rate field)
- The scheduling priority (Priority field)
- The maximum bandwidth in the priority group that the queue can consume (Shaping rate field)
- The drop profile loss priority (Loss priority field) for each drop profile name (name field)

The command output shows that:

- The scheduler map `be-map` has been created and has these properties:
 - There are two schedulers, `be-sched` and `nc-sched`.
 - The scheduler `be-sched` has one forwarding class, `best-effort`.
 - Scheduler `be-sched` forwarding class `best-effort` has a minimum guaranteed bandwidth of 90 percent, can consume a maximum of 100 percent of the priority group bandwidth, and uses the drop profile `dp-be-low` for low loss-priority traffic, the default drop profile for medium-high loss-priority traffic, and the drop profile `dp-be-high` for high loss-priority traffic.
 - The scheduler `nc-sched` has one forwarding class, `network-control`.
 - The `network-control` forwarding class has a minimum guaranteed bandwidth of 10 percent, can consume a maximum of 100 percent of the priority group bandwidth, and uses the drop profile `dp-nc` for low loss-priority traffic and the default drop profile for medium-high and high loss priority traffic.
- The scheduler map `n1-map` has been created and has these properties:
 - There are two schedulers, `fcoe-sched` and `n1-sched`.
 - The scheduler `fcoe-sched` has one forwarding class, `fcoe`.
 - The `fcoe` forwarding class has a minimum guaranteed bandwidth of 60 percent, and can consume a maximum of 100 percent of the priority group bandwidth.
 - The scheduler `n1-sched` has one forwarding class, `no-loss`.
 - The `no-loss` forwarding class has a minimum guaranteed bandwidth of 40 percent, and can consume a maximum of 100 percent of the priority group bandwidth.
- The scheduler map `mcast-map` has been created and has these properties:
 - There is one scheduler, `mcast-sched`.
 - The scheduler `mcast-sched` has one forwarding class, `mcast`.
 - The `mcast` forwarding class has a minimum guaranteed bandwidth of 100 percent, and can consume a maximum of 100 percent of the priority group bandwidth.

Verifying Fabric Queue Scheduler Configuration and Mapping on the Interconnect Device

Purpose

Verify that the Interconnect device ICD1 fabric queue schedulers are configured with the correct bandwidth parameters, mapped to the correct fabric `fc-sets`, and mapped to the correct drop profiles.

Action

List the fabric scheduler maps using the operational mode command `show class-of-service scheduler-map-forwarding-class-sets`:

```
user@switch> show class-of-service scheduler-map-forwarding-class-sets
```

```
Scheduler map forwarding class set: fab-traffic-map, Index: 2
```

```
Scheduler: fab-be-sched, Forwarding class set: fabric_fcset_be, Index: 21
```

```
Transmit rate: 25 percent, Rate Limit: none, Buffer size: 25 percent,
```

```
Buffer Limit: none, Priority: low
```

```
Excess Priority: unspecified
```

```
Shaping rate: 100 percent,
```

```
Drop profiles:
```

Loss priority	Protocol	Index	Name
Low	any	55387	fab-dp-be-low
Medium high	any	1	<default-drop-profile>
High	any	4369	fab-dp-be-high

```
Scheduler: fab-fcoe-sched, Forwarding class set: fabric_fcset_noloss1, Index: 23
```

```
Transmit rate: 30 percent, Rate Limit: none, Buffer size: 30 percent,
```

```
Buffer Limit: none, Priority: low
```

```
Excess Priority: unspecified
```

```
Shaping rate: 100 percent,
```

```
Drop profiles:
```

Loss priority	Protocol	Index	Name
Low	any	1	<default-drop-profile>
Medium high	any	1	<default-drop-profile>
High	any	1	<default-drop-profile>

```
Scheduler: fab-nl-sched, Forwarding class set: fabric_fcset_noloss2, Index: 27
```

```
Transmit rate: 25 percent, Rate Limit: none, Buffer size: 25 percent,
```

```
Buffer Limit: none, Priority: low
```

```
Excess Priority: unspecified
```

```
Shaping rate: 100 percent,
```

```
Drop profiles:
```

Loss priority	Protocol	Index	Name
Low	any	1	<default-drop-profile>
Medium high	any	1	<default-drop-profile>
High	any	1	<default-drop-profile>

```
Scheduler: fab-mcast-sched, Forwarding class set: fabric_fcset_multicast1, Index: 32
```



```

Transmit rate: 20 percent, Rate Limit: none, Buffer size: remainder,
Buffer Limit: none, Priority: low
Excess Priority: unspecified
Shaping rate: 100 percent,
Drop profiles:

```

Loss priority	Protocol	Index	Name
Low	any	1	<default-drop-profile>
Medium high	any	1	<default-drop-profile>
High	any	1	<default-drop-profile>

Meaning

The `show class-of-service scheduler-map-forwarding-class-sets` command lists the configured fabric scheduler map. The command output includes:

- The name of the fabric scheduler map (Scheduler map forwarding class set field)
- The name of the fabric scheduler (Scheduler field)
- The fabric fc-sets mapped to the scheduler (Forwarding class set field)
- The minimum guaranteed queue bandwidth (Transmit rate field)
- The maximum bandwidth in the priority group that the queue can consume (Shaping rate field)
- The drop profile loss priority (Loss priority field) for each drop profile name (Name field)

The command output shows that:

- The fabric scheduler map `fab-traffic-map` has been created and has these properties:
 - There are four fabric schedulers, `fab-be-sched`, `fab-fcoe-sched`, `fab-nl-sched`, and `fab-mcast-sched`.
 - The fabric scheduler `fab-be-sched` has one fabric fc-set, `fabric_fcset_be`.

The fabric fc-set `fabric_fcset_be` has a minimum guaranteed bandwidth of 25 percent, can consume a maximum of 100 percent of the priority group bandwidth, and uses the drop profile `fab-dp-be-low` for low loss-priority traffic, the default drop profile for medium-high loss-priority traffic, and the drop profile `fab-dp-be-high` for high loss-priority traffic.

- The fabric scheduler `fab-fcoe-sched` has one fabric fc-set, `fabric_fcset_noloss1`.

The `fabric_fcset_noloss1` fabric fc-set has a minimum guaranteed bandwidth of 30 percent, and can consume a maximum of 100 percent of the priority group bandwidth.

- The fabric scheduler `fab-nl-sched` has one fabric fc-set, `fabric_fcset_noloss2`.

The `fabric_fcset_noloss2` fabric fc-set has a minimum guaranteed bandwidth of 25 percent, and can consume a maximum of 100 percent of the priority group bandwidth.

- The fabric scheduler `fab-mcast-sched` has one fabric fc-set, `fabric_fcset_mcast1`.

The `fabric_fcset_multicast1` fabric fc-set has a minimum guaranteed bandwidth of 20 percent, and can consume a maximum of 100 percent of the priority group bandwidth.

Verifying Traffic Control Profile Configuration on the Node Devices

Purpose

Verify that the traffic control profiles (priority groups) `be-tcp`, `nl-tcp`, and `mcast-tcp` have been created with the correct bandwidth parameters and scheduler mapping.

Action

List the traffic control profiles using the operational mode command `show class-of-service traffic-control-profile`:

```
user@switch> show class-of-service traffic-control-profile
```

```
Traffic control profile: be-tcp, Index: 40535
```

```
  Shaping rate: 100 percent
```

```
  Scheduler map: be-map
```

```
  Guaranteed rate: 25 percent
```

```
Traffic control profile: nl-tcp, Index: 37959
```

```
  Shaping rate: 100 percent
```

```
  Scheduler map: nl-map
```

```
  Guaranteed rate: 50 percent
```

```
Traffic control profile: mcast-tcp, Index: 47661
```

```
  Shaping rate: 100 percent
```

```
  Scheduler map: mcast-map
```

```
  Guaranteed rate: 25 percent
```

Meaning

The `show class-of-service traffic-control-profile` command lists all of the configured traffic control profiles. For each traffic control profile, the command output includes:

- The name of the traffic control profile (Traffic control profile)

- The maximum port bandwidth the priority group can consume (Shaping rate)
- The scheduler map associated with the traffic control profile (Scheduler map)
- The minimum guaranteed priority group port bandwidth (Guaranteed rate)

The command output shows that:

- The traffic control profile `be-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `be-map`, and has a minimum guaranteed bandwidth of 25 percent of port bandwidth.
- The traffic control profile `nl-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `nl-map`, and has a minimum guaranteed bandwidth of 50 percent.
- The traffic control profile `mcast-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `mcast-map`, and has a minimum guaranteed bandwidth of 25 percent.

Verifying That PFC Is Enabled on Lossless Queues on the Node Devices

Purpose

Verify that PFC is enabled on the correct queues (as mapped to IEEE 802.1p priorities in the forwarding class configuration) for lossless transport.

Action

List the congestion notification profiles using the operational mode command `show class-of-service congestion-notification`:

```
user@switch> show class-of-service congestion-notification
Type: Input, Name: nl-cnp, Index: 51687
Priority    PFC
000        Disabled
001        Disabled
010        Disabled
011        Enabled
100        Enabled
101        Disabled
110        Disabled
111        Disabled
```

Meaning

The `show class-of-service congestion-notification` command lists all of the congestion notification profiles and the IEEE 802.1p code points with PFC enabled. The command output shows that PFC is enabled for code points 011 (fcoe queue) and 100 (no-loss queue) for the nl-cnp congestion notification profile.

Verifying Access and Fabric Interface Scheduling Configuration on the Node Devices

Purpose

Verify that the correct fc-sets, traffic control profiles, and congestion notification profiles are mapped to the correct interfaces on Node devices ND1 and ND2.

Action

List the interfaces on Node devices ND1 and ND2 using the operational mode command `show configuration class-of-service interfaces`. For example, the output on Node device ND1 shows:

```
user@switch> show configuration class-of-service interfaces
ND1:xe-0/0/20 {
    forwarding-class-set {
        best-effort-pg {
            output-traffic-control-profile be-tcp;
        }
        noloss-pg {
            output-traffic-control-profile nl-tcp;
        }
        multideestination-pg {
            output-traffic-control-profile mcast-tcp;
        }
    }
    congestion-notification-profile nl-cnp;
}
ND1:xe-0/0/21 {
    forwarding-class-set {
        best-effort-pg {
            output-traffic-control-profile be-tcp;
        }
        noloss-pg {
            output-traffic-control-profile nl-tcp;
        }
    }
}
```

```

        multidestination-pg {
            output-traffic-control-profile mcast-tcp;
        }
    }
    congestion-notification-profile nl-cnp;
}
ND1:fte-0/1/0 {
    forwarding-class-set {
        best-effort-pg {
            output-traffic-control-profile be-tcp;
        }
        noloss-pg {
            output-traffic-control-profile nl-tcp;
        }
        multidestination-pg {
            output-traffic-control-profile mcast-tcp;
        }
    }
}
}

```

Meaning

The `show configuration class-of-service interfaces` command shows that the fc-sets and (output) traffic control profiles mapped to the interfaces are:

- best-effort-pg fc-set with be-tcp traffic control profile
- noloss-pg fc-set with nl-tcp traffic control profile
- multidestination-pg fc-set with mcast-tcp traffic control profile

The command also shows that the access interfaces include the congestion notification profile nl-cnp to enable PFC on the IEEE 802.1p code points of lossless traffic.

Verifying Fabric Interface Scheduling Configuration on the Interconnect Device

Purpose

Verify that the correct fabric scheduler maps are associated with the correct fabric and Clos fabric interfaces on Interconnect device ICD1.

Action

List the interfaces using the operational mode command `show configuration class-of-service interfaces`:

```
user@switch> show configuration class-of-service interfaces
ICD1:fte-0/0/3 {
    scheduler-map-forwarding-class-set fab-traffic-map;
}
ICD1:fte-1/0/7
    scheduler-map-forwarding-class-set fab-traffic-map;
}
ICD1:bft-*/*/* {
    scheduler-map-forwarding-class-set fab-traffic-map;
}
```

Meaning

The `show configuration class-of-service interfaces` command shows that the same fabric forwarding class scheduler map is on all of the interfaces:

- `fab-traffic-map`

RELATED DOCUMENTATION

[Understanding CoS Scheduling Across the QFabric System | 327](#)

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

[Understanding CoS Fabric Forwarding Class Sets | 396](#)

[Understanding CoS Output Queue Schedulers | 186](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 223](#)

[Understanding Default CoS Scheduling and Classification | 95](#)

[Example: Configuring Queue Schedulers | 198](#)

[Example: Configuring WRED Drop Profiles | 300](#)

[Example: Configuring Drop Profile Maps | 307](#)

Understanding CoS Fabric Forwarding Class Sets

IN THIS SECTION

- [Default Fabric Forwarding Class Sets | 397](#)
- [Fabric Forwarding Class Set Configuration and Implementation | 401](#)
- [QFabric System CoS | 403](#)
- [Support for Flow Control and Lossless Transport Across the Fabric | 403](#)
- [Viewing Fabric Forwarding Class Set Information | 406](#)
- [Summary of Fabric Forwarding Class Set and Node Device Forwarding Class Set Differences | 407](#)

Fabric forwarding class sets (fabric fc-sets) are similar to the fc-sets (priority groups) you configure on Node devices. The major differences are:

1. Fabric fc-sets group traffic for transport across the QFX3008-I or QFX3600-I Interconnect device (the fabric). Node device fc-sets group traffic on a Node device for transport across that Node device.
2. Fabric fc-sets are global. They apply to the entire fabric. Node device fc-sets apply only to the Node device on which they are configured.
3. Fabric fc-sets are mapped directly to Interconnect device output queues; in this way, they behave similarly to forwarding classes on a Node device.

Fabric fc-sets map to Interconnect device fabric output queues statically—you cannot configure the mapping of fabric fc-sets to fabric output queues. All traffic in a fabric fc-set maps to the same output queue.

Node device fc-sets include forwarding classes that map to Node device output queues, and you can configure the mapping of forwarding classes to output queues (or you can use the default mapping). Because output queues are mapped to forwarding classes, different classes of traffic in a Node device fc-set can be mapped to different output queues.

Node device fc-sets consist of forwarding classes containing traffic that requires similar CoS treatment. (Forwarding classes are default forwarding classes or user-defined forwarding classes.) You can configure CoS for each fc-set to determine how the traffic of its forwarding classes is scheduled on a Node device.

When traffic exits a Node device interface and enters an Interconnect device fabric interface, the Interconnect device uses the same forwarding classes to group traffic. The forwarding classes are mapped to global fabric fc-sets for transport across the fabric. Like fc-sets on a Node device, fabric fc-

sets also contain traffic that requires similar CoS treatment. Also like fc-sets on a Node device, you can configure CoS on fabric fc-sets.

Fabric fc-sets reside on the Interconnect device and are global to the QFabric system. Fabric fc-sets apply to all traffic that traverses the fabric. The mapping of forwarding classes to fabric fc-sets is global and applies to all forwarding classes with traffic that traverses the fabric from all connected Node devices. You can change the mapping of forwarding classes to fabric fc-sets. All mapping changes you make are global. For example, if you change the fabric fc-set to forwarding class mapping of the default best-effort forwarding class, then every Node device's best-effort forwarding class traffic that traverses the fabric is mapped to that fabric fc-set. The CoS you configure on a fabric fc-set applies to all the traffic that belongs to that fabric fc-set, from all connected Node devices.

This topic describes:

Default Fabric Forwarding Class Sets

Interconnect devices have 12 default fabric fc-sets, including five visible default fabric fc-sets, four for unicast traffic and one for multideestination (multicast, broadcast, and destination lookup failure) traffic.

There are also seven hidden default fabric fc-sets. There are three hidden default fabric fc-sets for multideestination traffic that you can use if you want to map different multideestination forwarding classes to different multideestination fabric fc-sets. There are four hidden default fabric fc-sets for lossless traffic that you can use to map different lossless forwarding classes (priorities) to different lossless fabric fc-sets. [Table 78 on page 397](#) shows the default fabric fc-sets:

Table 78: Default Fabric Forwarding Class Sets

Fabric Forwarding Class Set Name	Characteristics
fabric_fcset_be	Transports best-effort unicast traffic across the fabric.
fabric_fcset_strict_high	Transports unicast traffic that has been configured with strict-high priority and in the network-control forwarding class across the fabric. This fabric fc-set receives as much bandwidth across the fabric as it needs to service the traffic in the group up to the entire fabric interface bandwidth. For this reason, exercise caution when mapping traffic to this fabric fc-set to avoid starving other traffic.
fabric_fcset_noloss1	Transports unicast traffic in the default fcoe forwarding class across the fabric.

Table 78: Default Fabric Forwarding Class Sets (Continued)

Fabric Forwarding Class Set Name	Characteristics
fabric_fcset_noloss2	Transports unicast traffic in the default no-loss forwarding class across the fabric.
fabric_fcset_noloss3	(Hidden) No traffic is assigned by default to this fabric fc-set. Unless traffic is mapped to this fabric fc-set, this fabric fc-set remains hidden. This fabric fc-set is valid only for lossless forwarding classes.
fabric_fcset_noloss4	(Hidden) No traffic is assigned by default to this fabric fc-set. Unless traffic is mapped to this fabric fc-set, this fabric fc-set remains hidden. This fabric fc-set is valid only for lossless forwarding classes.
fabric_fcset_noloss5	(Hidden) No traffic is assigned by default to this fabric fc-set. Unless traffic is mapped to this fabric fc-set, this fabric fc-set remains hidden. This fabric fc-set is valid only for lossless forwarding classes.
fabric_fcset_noloss6	(Hidden) No traffic is assigned by default to this fabric fc-set. Unless traffic is mapped to this fabric fc-set, this fabric fc-set remains hidden. This fabric fc-set is valid only for lossless forwarding classes.
fabric_fcset_multicast1	Transports multdestination traffic in the mcast forwarding class across the fabric. This fabric fc-set is valid only for multdestination forwarding classes.
fabric_fcset_multicast2	(Hidden) No traffic is assigned by default to this fabric fc-set. Unless traffic is mapped to this fabric fc-set, this fabric fc-set remains hidden. This fabric fc-set is valid only for multdestination forwarding classes.
fabric_fcset_multicast3	(Hidden) No traffic is assigned by default to this fabric fc-set. Unless traffic is mapped to this fabric fc-set, this fabric fc-set remains hidden. This fabric fc-set is valid only for multdestination forwarding classes.

Table 78: Default Fabric Forwarding Class Sets (Continued)

Fabric Forwarding Class Set Name	Characteristics
fabric_fcset_multicast4	(Hidden) No traffic is assigned by default to this fabric fc-set. Unless traffic is mapped to this fabric fc-set, this fabric fc-set remains hidden. This fabric fc-set is valid only for multideestination forwarding classes.

The five default forwarding classes (best-effort, fcoe, no-loss, network-control, and mcast) are mapped to the fabric fc-sets by default as shown in [Table 79 on page 399](#).

Table 79: Default Forwarding Class to Fabric Forwarding Class Set Mapping

Forwarding Class	Fabric Forwarding Class Set	Fabric Output Queue	Maximum MTU Supported for Lossless Operation
best-effort	fabric_fcset_be	0	NA
network-control	fabric_fcset_strict_high	7	NA
fcoe	fabric_fcset_noloss1	1	9K
no-loss	fabric_fcset_noloss2	2	9K
mcast	fabric_fcset_multicast1	8	NA
No forwarding classes are mapped by default to this hidden fabric fc-set.	fabric_fcset_noloss3	3	9k
No forwarding classes are mapped by default to this hidden fabric fc-set.	fabric_fcset_noloss4	4	9k
No forwarding classes are mapped by default to this hidden fabric fc-set.	fabric_fcset_noloss5	5	9k

Table 79: Default Forwarding Class to Fabric Forwarding Class Set Mapping (Continued)

Forwarding Class	Fabric Forwarding Class Set	Fabric Output Queue	Maximum MTU Supported for Lossless Operation
No forwarding classes are mapped by default to this hidden fabric fc-set.	fabric_fcset_noloss6	6	9k
No forwarding classes are mapped by default to this hidden fabric fc-set.	fabric_fcset_multicast2	9	NA
No forwarding classes are mapped by default to this hidden fabric fc-set.	fabric_fcset_multicast3	10	NA
No forwarding classes are mapped by default to this hidden fabric fc-set.	fabric_fcset_multicast4	11	NA

The maximum fiber cable length between the QFabric system Node device and the QFabric system Interconnect device is 150 meters.

TIP: If you explicitly configure lossless forwarding classes, we recommend that you map each user-configured lossless forwarding class to an unused fabric fc-set (fabric_fcset_noloss3 through fabric_fcset_noloss6) on a one-to-one basis: one lossless forwarding class mapped to one lossless fabric fc-set.

The reason for one-to-one mapping is to avoid fate sharing of lossless flows. Because each fabric fc-set is mapped statically to an output queue, when you map more than one forwarding class to a fabric fc-set, all of the traffic in all of the forwarding classes that belong to the fabric fc-set uses the same output queue. If that output queue becomes congested due to congestion caused by one of the flows, the other flows are also affected. (They share fate because the flow that congests the output queue affects flows that are not experiencing congestion.)

If you want to map different multdestination forwarding classes to different multdestination fabric fc-sets, use one or more of the hidden multdestination fabric fc-sets.

NOTE: The global mapping of forwarding classes to fabric fc-sets is independent of the mapping of forwarding classes to Node device fc-sets. Global mapping of forwarding classes to fabric fc-sets occurs only on the Interconnect device. The Node device mapping of forwarding classes to

fc-sets does not affect the global mapping of forwarding classes to fabric fc-sets on the Interconnect device, and vice versa.

When you define new forwarding classes on a Node device, you explicitly map those forwarding classes to Node device fc-sets. However, new (user-created) forwarding classes are mapped by default to fabric fc-sets. (You can override the default mapping if you want to configure the forwarding class to fabric fc-set mapping explicitly, as described in the next section.)

By default:

- All best-effort traffic forwarding classes that you create are mapped to the `fabric_fcset_be` fabric fc-set.
- All lossless traffic forwarding classes that you create are mapped to the `fabric_fcset_noloss1` or `fabric_fcset_noloss2` fabric fc-set.

NOTE: To avoid fate sharing, we recommend that you configure one-to-one mapping of user-configured lossless forwarding classes to lossless fabric fc-sets instead of using the default mapping. You can also use firewall filters to mitigate fate sharing by separating flows that belong to the same forwarding class as the traffic traverses the Interconnect device (see [Understanding How to Mitigate Fate Sharing on a QFabric System Interconnect Device by Remapping Traffic Flows \(Forwarding Classes\)](#) for more information.)

- All multidestination traffic forwarding classes that you create are mapped to the `fabric_fcset_multicast1` fabric fc-set.
- All strict-high priority traffic and network-control forwarding classes that you create are mapped to the `fabric_fcset_strict_high` fabric fc-set.

Fabric Forwarding Class Set Configuration and Implementation

You can map forwarding classes to fabric fc-sets and configure CoS scheduling for fabric fc-sets. This section describes:

Mapping Forwarding Classes to Fabric Forwarding Class Sets

If you do not want to use the default mapping of forwarding classes to fabric fc-sets, you can map forwarding classes to fabric fc-sets in the same way as you map forwarding classes to Node device fc-sets. To do this, use exactly the same statement that you use to map forwarding classes to fc-sets, but instead of specifying a Node device fc-set name, specify a fabric fc-set name.

NOTE: The global mapping of forwarding classes to fabric fc-sets does not affect the mapping of forwarding classes to Node device fc-sets. The global forwarding class mapping to fabric fc-sets pertains to the traffic only when it enters, traverses, and exits the fabric. The forwarding class mapping to fc-sets on a Node device is valid within that Node device.

Mapping forwarding classes to fabric fc-sets does not affect the scheduling configuration of the forwarding classes or fc-sets on Node devices. Fabric fc-set scheduling pertains to traffic only when it enters, traverses, and exits the Interconnect device fabric.

If you change the mapping of a forwarding class to a fabric fc-set, the new mapping is global and applies to all traffic in that forwarding class, regardless of which Node device forwards the traffic to the Interconnect device.

- To assign one or more forwarding classes to a fabric fc-set:

```
[edit class-of-service]
user@switch# set forwarding-class-sets fabric-forwarding-class-set-name class forwarding-  
class-name
```

For example, to map a user-defined forwarding class named best-effort-2 to the fabric fc-set fabric_fcset_be:

```
[edit class-of-service]
user@switch# set forwarding-class-sets fabric_fcset_be class best-effort-2
```

NOTE: Because fabric fc-set configuration is global, in this example all forwarding classes with the name best-effort-2 on all of the Node devices connected to the fabric use the fabric_fcset_be fabric fc-set to transport traffic across the fabric.

Fabric Forwarding Class Set Implementation

The following rules apply to fabric fc-sets:

- You cannot create new fabric fc-sets. Only the twelve default fabric fc-sets are available.
- You cannot delete a default fabric fc-set.
- You cannot attach a fabric fc-set to a Node device interface. Fabric fc-sets are used only on the Interconnect device fabric, not on Node devices.

- You can map only multdestination forwarding classes to multdestination fabric fc-sets.
- You cannot map multdestination forwarding classes to unicast fabric fc-sets.
- You cannot map unicast forwarding classes to multdestination fabric fc-sets.

QFabric System CoS

When traffic enters and exits the same QFabric system Node device, CoS works the same as it works on a standalone switch.

However, when traffic enters a QFabric system Node device, crosses the Interconnect device, and then exits a different Node device, CoS is applied differently:

1. Traffic entering the ingress Node device receives the CoS configured at the Node ingress (packet classification and congestion notification profile for PFC).
2. When traffic goes from the ingress Node device to the Interconnect device, the fabric fc-set CoS is applied to the traffic.
3. When traffic goes from the Interconnect device to the egress Node device, the egress Node device applies CoS at the egress port (egress queue scheduling, WRED, and IEEE 802.1p or DSCP code-point rewrite).

Traffic that traverses the Interconnect device can use the default CoS fabric scheduling or you can configure two-tier hierarchical CoS scheduling explicitly on fabric fc-sets as described in [Understanding CoS Scheduling Across the QFabric System](#).

Support for Flow Control and Lossless Transport Across the Fabric

The Interconnect device incorporates flow control mechanisms to support lossless transport during periods of congestion on the fabric. To support the priority-based flow control (PFC) feature on the Node devices, the fabric interfaces use LLFC to support lossless transport for up to six IEEE 802.1p priorities when the following two configuration constraints are met:

1. The IEEE 802.1p priority used for the traffic that requires lossless transport is mapped to a lossless forwarding class on the Node devices.
2. The lossless forwarding class must be mapped to a lossless fabric fc-set on the Interconnect device (fabric_fcset_noloss1, fabric_fcset_noloss2, fabric_fcset_noloss3, fabric_fcset_noloss4, fabric_fcset_noloss5, or fabric_fcset_noloss6).

When traffic meets the two configuration constraints, the fabric propagates the back pressure from the egress Node device across the fabric to the ingress Node device during periods of congestion. However, to achieve end-to-end lossless transport across the switch, you must also configure a congestion notification profile to enable PFC on the Node device ingress ports.

For all other combinations of IEEE 802.1p priority to forwarding class mapping and all other combinations of forwarding class to fabric fc-set mapping, the congestion control mechanism is normal packet drop. For example:

- **Case 1**—If the IEEE 802.1p priority 5 is mapped to the lossless fcoe forwarding class, and the fcoe forwarding class is mapped to the fabric_fcset_noloss1 fabric fc-set, then the congestion control mechanism is PFC.
- **Case 2**—If the IEEE 802.1p priority 5 is mapped to the lossless fcoe forwarding class, and the fcoe forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop.
- **Case 3**—If the IEEE 802.1p priority 5 is mapped to the lossless no-loss forwarding class, and the no-loss forwarding class is mapped to the fabric_fcset_noloss2 fabric fc-set, then the congestion control mechanism is PFC.
- **Case 4**—If the IEEE 802.1p priority 5 is mapped to the lossless no-loss forwarding class, and the no-loss forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop.
- **Case 5**—If the IEEE 802.1p priority 5 is mapped to the best-effort forwarding class, and the best-effort forwarding class is mapped to the fabric_fcset_be fabric fc-set, then the congestion control mechanism is packet drop.
- **Case 6**—If the IEEE 802.1p priority 5 is mapped to the best-effort forwarding class, and the best-effort forwarding class is mapped to the fabric_fcset_noloss1 fabric fc-set, then the congestion control mechanism is packet drop.

NOTE: Lossless transport across the fabric also must meet the following two conditions:

1. The maximum cable length between the Node device and the Interconnect device is a 150 meters of fiber cable.
2. The maximum frame size is 9216 bytes.

If the MTU is 9216 KB, in some cases the QFabric system supports only five lossless forwarding classes instead of six lossless forwarding classes because of headroom buffer limitations.

The number of IEEE 802.1p priorities (forwarding classes) the QFabric system can support for lossless transport across the Interconnect device fabric depends on several factors:

- **Approximate fiber cable length**—The longer the fiber cable that connects Node device fabric (FTE) ports to the Interconnect device fabric ports, the more data the connected ports need to buffer when a pause is asserted. (The longer the fiber cable, the more frames are traversing the cable when

a pause is asserted. Each port must be able to store all of the “in transit” frames in the buffer to preserve lossless behavior and avoid dropping frames.)

- MTU size—The larger the maximum frame sizes the buffer must hold, the fewer frames the buffer can hold. The larger the MTU size, the more buffer space each frame consumes.
- Total number of Node device fabric ports connected to the Interconnect device—The higher the number of connected fabric ports, the more headroom buffer space the Node device needs on those fabric ports to support the lossless flows that traverse the Interconnect device. Because more buffer space is used on the Node device fabric ports, less buffer space is available for the Node device access ports, and a lower total number of lossless flows are supported.

The QFabric system supports six lossless priorities (forwarding classes) under most conditions. The priority group headroom that remains after allocating headroom to lossless flows is sufficient to support best-effort and multidestination traffic.

[Table 80 on page 405](#) shows how many lossless priorities the QFabric system supports under different conditions (fiber cable lengths and MTUs) in cases when the QFabric system supports fewer than six lossless priorities. The number of lossless priorities is the same regardless of how many Node device FTE ports are connected to the Interconnect device. However, the higher the number of FTE ports connected to the Interconnect device, the lower the number of total lossless flows supported. In all cases that are not shown in [Table 80 on page 405](#), the QFabric system supports six lossless priorities.

NOTE: The system does not perform a configuration commit check that compares available system resources with the number of lossless forwarding classes configured. If you commit a configuration with more lossless forwarding classes than the system resources can support, frames in lossless forwarding classes might be dropped.

Table 80: Lossless Priority (Forwarding Class) Support for Node Devices When Fewer than Six Lossless Priorities Are Supported

MTU in Bytes	Fiber Cable Length in Meters (Approximate)	Maximum Number of Lossless Priorities (Forwarding Classes) on the Node Device
9216 (9K)	100	5
9216 (9K)	150	5

NOTE: The total number of lossless flows decreases as resource consumption increases. For a Node device, the higher the number of FTE ports connected to the Interconnect device, the larger the MTU, and the longer the fiber cable length, the fewer total lossless flows the QFabric system can support.

Viewing Fabric Forwarding Class Set Information

You can display information about fabric fc-sets using the same CLI command you use to display information about Node device fc-sets:

```
user@switch> show class-of-service forwarding-class-set
```

Forwarding class set: fabric_fcset_be, Type: fabric-type, Forwarding class set index: 1	
Forwarding class	Index
best-effort	0

Forwarding class set: fabric_fcset_mcast1, Type: fabric-type, Forwarding class set index: 5	
Forwarding class	Index
mcast	8

Forwarding class set: fabric_fcset_mcast2, Type: fabric-type, Forwarding class set index: 6	
---	--

Forwarding class set: fabric_fcset_mcast3, Type: fabric-type, Forwarding class set index: 7	
---	--

Forwarding class set: fabric_fcset_mcast4, Type: fabric-type, Forwarding class set index: 8	
---	--

Forwarding class set: fabric_fcset_noloss1, Type: fabric-type, Forwarding class set index: 2	
Forwarding class	Index
fcoe	1

Forwarding class set: fabric_fcset_noloss2, Type: fabric-type, Forwarding class set index: 3	
Forwarding class	Index
no-loss	2

Forwarding class set: fabric_fcset_noloss3, Type: fabric-type, Forwarding class set index: 9	
--	--

Forwarding class set: fabric_fcset_noloss4, Type: fabric-type, Forwarding class set index: 10	
---	--

Forwarding class set: fabric_fcset_noloss5, Type: fabric-type, Forwarding class set index: 11	
---	--

Forwarding class set: fabric_fcset_noloss6, Type: fabric-type, Forwarding class set index: 12

Forwarding class set: fabric_fcset_strict_high, Type: fabric-type, Forwarding class set index: 4

Forwarding class	Index
network-control	3

[Table 81 on page 407](#) describes the meaning of the show class-of-service forwarding-class-set output fields when you display fabric fc-set information.

Table 81: show class-of-service forwarding-class-set Command Output Fields

Field Name	Field Description
Forwarding class set	Name of the fabric forwarding class set.
Type	Type of forwarding class set: <ul style="list-style-type: none"> Fabric-type—Fabric fc-set Normal-type—Node device fc-set
Forwarding class set index	Index of this forwarding class set.
Forwarding class	Name of a forwarding class.
Index	Index of the forwarding class.

Summary of Fabric Forwarding Class Set and Node Device Forwarding Class Set Differences

[Table 82 on page 407](#) summarizes the differences between fabric fc-sets and fc-sets:

Table 82: Summary of Differences Between Fabric fc-sets and Local fc-sets

Characteristic	Fabric fc-set	Local fc-set
Location	QFX3008-I or QFX3600-I Interconnect device (the fabric).	QFabric Node device.

Table 82: Summary of Differences Between Fabric fc-sets and Local fc-sets (Continued)

Characteristic	Fabric fc-set	Local fc-set
Global or local	Global, valid for the entire fabric.	Local to the Node device on which the fc-set is configured.
Ability to create (define) a new fc-set	No. Use the 12 default fabric fc-sets provided.	Yes.
Ability to configure CoS	User-configurable using fabric fc-set scheduler maps.	User-configurable using traffic control profiles.
Ability to map forwarding classes to an fc-set	Yes. Mapping is global and applies to all forwarding classes across the Interconnect device fabric (traffic from all connected Node devices).	Yes. Mapping is local to a Node device and applies only to the forwarding classes on the Node device.

RELATED DOCUMENTATION
[Understanding CoS Forwarding Class Sets \(Priority Groups\)](#)
[Understanding CoS Scheduling Across the QFabric System](#)
[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports](#)
[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\)](#)
[Understanding How to Mitigate Fate Sharing on a QFabric System Interconnect Device by Remapping Traffic Flows \(Forwarding Classes\)](#)
[Defining CoS Forwarding Class Sets](#)
[Example: Configuring Forwarding Class Sets](#)
[Example: Configuring CoS Scheduling Across the QFabric System](#)
[show class-of-service forwarding-class-set](#)

Configuring CoS Fabric Forwarding Class Set Scheduler Maps (Fabric Scheduler to Fabric FC-Set Mapping)

Fabric forwarding class set scheduler maps map fabric fc-sets to fabric schedulers, much the same way that scheduler maps on Node devices map fc-sets to schedulers. Each fabric fc-set represents an output queue on the Interconnect device interfaces, and each fabric forwarding class set scheduler sets bandwidth scheduling and other CoS properties.

When you map a fabric fc-set to a fabric scheduler, the fabric scheduler CoS properties determine the port bandwidth resources that the fabric fc-set traffic receives. You associate the fabric forwarding class set scheduler map a fabric interface to apply the scheduling properties to the traffic in the fabric fc-set on that fabric interface.

Using different fabric scheduler forwarding class set maps, you can map different schedulers to the same traffic (the same fabric fc-set) so that you can apply different scheduling to the traffic on different interfaces.

You can associate one fabric forwarding class set scheduler map to a fabric interface. Each fabric forwarding class scheduler map can contain multiple schedulers mapped to multiple fabric fc-sets.

Before you begin, you need to configure one or more fabric schedulers to map to fabric fc-sets. If you do not configure fabric schedulers, the system uses the default fabric schedulers and the default fabric forwarding class set scheduler mapping to fabric fc-sets.

To configure a fabric forwarding class set scheduler map using the CLI, define a name for the fabric forwarding class set scheduler map, and specify the fabric fc-set and fabric scheduler that you are mapping:

- ```
[edit class-of-service]
user@switch# set scheduler-map-forwarding-class-sets fabric-scheduler-map-name forwarding-
class-set fabric-fc-set-name scheduler fabric scheduler-name
```

For example, to configure a fabric scheduler map named `fab-be-map` that has a fabric fc-set named `fabric_fcset_be` mapped to fabric scheduler `fab-be-sched`:

```
[edit class-of-service]
user@switch# set scheduler-map-forwarding-class-sets fab-be-map forwarding-class-set
fabric_fcset_be scheduler fab-be-sched
```

## RELATED DOCUMENTATION

*Example: Configuring Queue Schedulers*

[Example: Configuring CoS Scheduling Across the QFabric System | 353](#)

[Understanding CoS Scheduling Across the QFabric System | 327](#)

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

*Understanding CoS Fabric Forwarding Class Sets*

## Understanding How to Mitigate Fate Sharing on a QFabric System Interconnect Device by Remapping Traffic Flows (Forwarding Classes)

### IN THIS SECTION

- [Fate Sharing on the Interconnect Device | 411](#)
- [Scenario 1: How Fate Sharing Can Occur on a QFabric System Interconnect Device | 414](#)
- [Scenario 2: How Forwarding Class Remapping Mitigates Fate Sharing on a QFabric System Interconnect Device | 416](#)
- [Fate Sharing Mitigation Process | 418](#)
- [Best Practices | 430](#)
- [Limitations and Notes on Behavior | 431](#)

On a QFabric system, traffic either is switched locally on a Node device (traffic enters and exits the same Node device without crossing the Interconnect device), or is switched remotely, across the Interconnect device. Traffic flows that belong to the same forwarding class are mapped to the same output queue and share the output queue resources. If congestion occurs on one of these flows, the congestion can affect the uncongested flows in the forwarding class when the flows use the same ingress interface.

For example, if a congested flow is paused to prevent packet loss, uncongested flows that use the same ingress interface are also paused because they share the same forwarding class and output queue. When a congested flow affects an uncongested flow, the flows share the same fate—this is known as fate sharing.

Fate sharing happens because pausing traffic is based on forwarding class. When a flow experiences congestion, the output queue sends a pause message to the input queue on which the flow arrived. On

that input queue, the pause message affects all traffic in the forwarding class that is mapped to the congested output queue. So all traffic in that forwarding class is paused on the input queue, not just the flow that is experiencing the congestion. This is how uncongested flows can share fate with a congested flow.

Traffic from many QFabric system Node devices crosses the Interconnect device, so flows within a given forwarding class are aggregated on the Interconnect device. The aggregated flows use the same output queue on the Interconnect device and are subject to fate sharing if the flows also use the same ingress interface.

In addition to the external physical interfaces that connect the Interconnect device to Node devices, the Interconnect device has internal Clos interfaces. The Interconnect device automatically selects the best path through its internal Clos interfaces. Path selection through the internal Clos interfaces is not configurable, so you cannot control the traffic that enters any particular ingress Clos interface, and so fate sharing can occur on the Interconnect device. (On Node devices, you control the traffic connected to an ingress interface, but on the Interconnect device, you cannot control which flows use a particular internal ingress Clos interface.)

However, to mitigate fate sharing on the Interconnect device, you can use firewall filters to separate the traffic in one forwarding class and split it into different forwarding classes. Remapping the flows into different forwarding classes means the flows use different output queues on the Interconnect device. If the flows use the same ingress interface on the Interconnect device, they do not experience fate sharing because only the flows mapped to the congested queue are paused, while the flows remapped to other forwarding classes (remapped to different queues) are not paused.

Mitigating fate sharing is often useful for lossless flows such as storage traffic, but is not limited to lossless flows. You can remap forwarding classes to mitigate fate sharing on the Interconnect device to separate flows that belong to any application (for example, iSCSI, NAS, FCoE, and so on), even when the flows are in the same VLAN.

## Fate Sharing on the Interconnect Device

Fate sharing can occur when multiple flows use the same output queue (the flows are mapped to the same forwarding class) and the same ingress interface. If the flows use the same ingress interface, then if one congested flow is paused, the uncongested flows in the same forwarding class are also paused at the shared ingress interface—the uncongested flows share fate with the congested flow. On an Interconnect device, the flows from many Node devices are aggregated, so the number of flows assigned to a given forwarding class and forwarded through a particular egress interface can be much greater on an Interconnect device than on a single Node device.

**NOTE:** The possibility of fate sharing cannot be avoided on Node device ingress interfaces. If two servers access a Node device on the same ingress interface, and both servers send traffic

flows that are classified into the same forwarding class (for example, if both flows are FCoE traffic and are classified into the *fcoe* forwarding class), then even if the flows are in different VLANs, congestion on one flow affects the other flow. The congested flow affects the uncongested flow because both flows share the same forwarding class (and therefore the same output queue and IEEE 802.1p code point), and priority-based flow control (PFC) is applied to the ingress interface, not to the VLAN. So when PFC pauses the congested flow on the ingress interface, the uncongested flow that uses the same code point is also paused.

An example of fate sharing is when two Fibre Channel over Ethernet (FCoE) flows are in the same forwarding class (so they use the same output queue) and use the same Interconnect device ingress interface. If one of those flows experiences congestion and the other flow does not experience congestion, the congested flow can affect the uncongested flow. If backpressure forces the ingress interface to pause the congested FCoE flow, the uncongested FCoE flow is also paused because the two flows use the same forwarding class (output queue) and traffic in that forwarding class is paused on the ingress interface.

Remapping flows that belong to the same forwarding class into different forwarding classes for transport across the Interconnect device mitigates fate sharing by separating the flows onto different output queues. Using different output queues means that the flows use different forwarding classes on the ingress interface. When a flow on one queue is paused, it does not affect flows that have been remapped onto other queues. The congestion only affects traffic on the paused queue, so the congestion only affects the congested forwarding class on the ingress interface.

After the traffic crosses the Interconnect device, the Node device on which the traffic egresses the QFabric system must map the traffic back to its original forwarding class before forwarding the traffic toward its destination, because the original forwarding class contains similar traffic, and is classified to support the CoS that the traffic type requires and the destination device expects.

For example, traffic destined for different targets in the same storage area network (SAN) normally should be in the same forwarding class, because a SAN uses one IEEE 802.1p code point (priority) to identify all traffic of a particular type, such as FCoE traffic. So when traffic destined for the SAN leaves the QFabric system, all of it must be mapped to the same forwarding class so that it uses the same IEEE 802.1p code point and is identified and classified the same way when it enters the SAN. This is why the QFabric system must map the traffic back into its original forwarding class after the traffic crosses the Interconnect device.

The QFabric system uses a firewall filter to remap traffic to a different forwarding class before it crosses the Interconnect device, and then map traffic back to its original forwarding class after it exits the Interconnect device.

The firewall filter must remap forwarding classes in each direction the traffic flows through the Interconnect device. For example, filter terms must remap traffic when it travels from a server to a target, and also when traffic travels from a target to a server. For each direction of traffic, you configure

a filter term that maps traffic into a different forwarding class when it enters the Interconnect device, and a filter term that maps traffic back into its original forwarding class after it exits the Interconnect device.

As with all firewall filters, there is a default discard at the end of each filter rule, so if you do not want to discard all traffic that is not explicitly permitted, you should add a final term to accept traffic that is not affected by the other terms. This is especially important when you are not remapping all of the traffic in a VLAN.

The QFabric system supports up to six lossless traffic classes called fabric forwarding class sets (fabric fc-sets) on the Interconnect device. (You can also configure up to six lossless forwarding classes on a Node device.) Each fabric fc-set maps to a different output queue on the Interconnect device.

**NOTE:** Fabric fc-sets on the Interconnect device are analogous to forwarding classes on Node devices, in that both fabric fc-sets and forwarding classes map to output queues on their respective devices. However, more than one forwarding class can map to a fabric fc-set, so a fabric fc-set can aggregate forwarding classes on the Interconnect device.

Fabric fc-set names are not user-configurable, and you cannot configure new fabric fc-sets. You can configure forwarding class to fabric fc-set mapping, so each fabric fc-set transports the traffic that you want it to transport.

The six lossless fabric fc-sets enable you to separate traffic from Node devices into as many as six lossless traffic classes on the Interconnect device. Each fabric fc-set uses a different output queue, so the flows (forwarding classes) mapped to each fabric fc-set use different ingress interfaces, and the flows in one fabric fc-set do not share fate with flows in other fabric fc-sets.

In addition, there are four multidestination fabric fc-sets on the Interconnect to handle multicast, broadcast, and destination lookup fail traffic.

**NOTE:** The flows (forwarding classes) within a fabric fc-set can share fate if they use the same ingress interface (the shared ingress interface could be an external 40-Gigabit interface or an internal Clos interface) because they use the same output queue, so they belong to the same forwarding class. However, the ability to separate flows into different forwarding classes enables you to spread the traffic among multiple output queues, and thus to mitigate the possibility of fate sharing because only traffic that belongs to the paused forwarding class (output queue) is paused on the ingress interface. Traffic remapped into other forwarding classes is not paused.



## Scenario 1: How Fate Sharing Can Occur on a QFabric System Interconnect Device

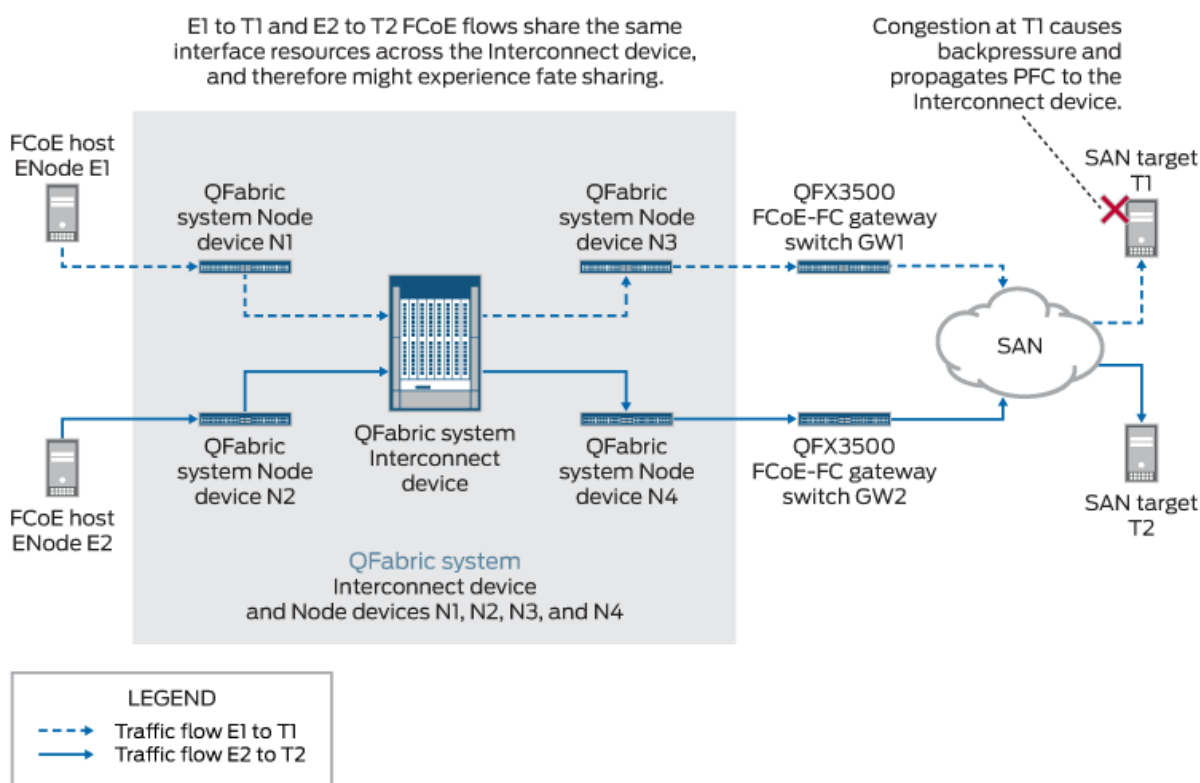
An example of traffic that might share fate across the Interconnect device is storage traffic. This scenario uses FCoE as an example.

**NOTE:** Any type of traffic that shares the same forwarding class (output queue) and Interconnect device ingress interface can experience fate sharing.

QFabric system Node devices aggregate FCoE traffic from connected ENodes. Because the FCoE traffic requires the same treatment across the network, in this scenario the FCoE traffic uses the same forwarding class on all of the Node devices (the default `fcoe` forwarding class), and is mapped to the same output queue on all of the Node devices. Because the Fibre Channel (FC) SAN usually expects traffic to have a priority value of 3 (IEEE 802.1p code point 011), priority 3 identifies all of the FCoE traffic.

All of the FCoE traffic that is not locally switched on the Node devices is remotely switched across the Interconnect device. A large amount of FCoE traffic might be switched across the Interconnect device, and all of that traffic uses the same egress queue. Whenever the FCoE flows use the same Interconnect device ingress interface, fate sharing can occur, as shown in [Figure 16 on page 415](#).

Figure 16: Fate Sharing Scenario: FCoE Traffic Shares Fate on the Interconnect Device



In [Figure 16 on page 415](#), FCoE traffic flows from two FCoE hosts (ENode E1 and ENode E2) through the QFabric system and an FCoE-FC gateway switch to two storage target devices in the SAN (target T1 and target T2). Target device T1 is experiencing congestion. Target device T2 is not experiencing congestion, as shown by the red “X”.

The dotted line shows the path that FCoE traffic from ENode E1 takes, entering the QFabric system at ingress Node device N1, flowing through the Interconnect device to the egress Node device N3, exiting the QFabric system to FCoE-FC gateway switch GW1, entering the SAN, then finally reaching target T1.

The solid line shows the path that FCoE traffic from ENode E2 takes, entering the QFabric system at ingress Node device N2, flowing through the Interconnect device to the egress Node device N4, exiting the QFabric system to FCoE-FC gateway switch GW2, entering the SAN, then finally reaching target T2.

When FCoE traffic from hosts ENode E1 and ENode E2 crosses the Interconnect device, the flows from the two hosts might use the same ingress interface at any of the Interconnect device interface stages (external 40-Gigabit interfaces or internal Clos interfaces). If the flows use the same ingress interface at any point, the paths of the flows converge at that interface on the input queue instead of remaining separate. (The dotted line and the solid line can intersect if at any time they share a common Interconnect device ingress interface.) When traffic flows assigned to the same forwarding class use the same ingress interface, fate sharing can occur because the flows use the same output queue.

In this scenario, the flows from hosts E1 and E2 share an ingress interface as they cross the Interconnect device. When target T1 experiences congestion, it sends a pause message to temporarily stop the incoming flow until the congestion clears, in order to prevent packet loss due to queue overfill. The pause message propagates back through the data path. Eventually, host E1 will receive a pause message and temporarily stop transmitting.

However, when the pause message reaches the Interconnect device ingress interface that the FCoE flows from hosts E1 and E2 share, not only is the flow originating from host E1 paused, the flow originating from host E2 is also paused, even though the E2 host flow is not experiencing congestion. Both flows are paused because both flows belong to the same forwarding class, and therefore use the same output queue, and use the same ingress interface. When the message pauses the E1-to-T1 flow, it also pauses the E2-to-T2 flow, because the all flows in the forwarding class are paused on the shared ingress interface, regardless of whether or not an individual flow in that forwarding class is experiencing congestion.

In this scenario, the uncongested FCoE flow from E2-to-T2 shares the same fate as the congested FCoE flow from E1-to-T1.

**NOTE:** This FCoE traffic scenario is one example of fate sharing. Fate sharing can occur on any flows that are mapped to the same forwarding class (output queue) and use the same Interconnect device ingress interface.

## Scenario 2: How Forwarding Class Remapping Mitigates Fate Sharing on a QFabric System Interconnect Device

Fate sharing occurs when traffic flows are assigned to the same forwarding class (and therefore use the same output queue) and also use the same ingress interface. The trick to mitigating the effects of fate sharing is to do one of two things: either ensure that flows assigned to the same forwarding class use different ingress interfaces, or ensure that flows use different forwarding classes, so that if they use the same ingress interface, they use different egress queues on the interface.

Ensuring that flows assigned to the same forwarding class use different ingress interfaces is not possible because the Interconnect device automatically selects the best path through its internal Clos interfaces. You cannot configure the Interconnect device to route traffic along a particular path within the device. However, you can separate the traffic assigned to one forwarding class into multiple forwarding classes for the journey across the Interconnect device. Remapping the flows into different forwarding classes means the flows use different output queues, so if the flows use the same ingress interface, they will not experience fate sharing.

**Figure 17: Fate Sharing Mitigation Scenario: FCoE Traffic Avoids Fate Sharing on the Interconnect Device**

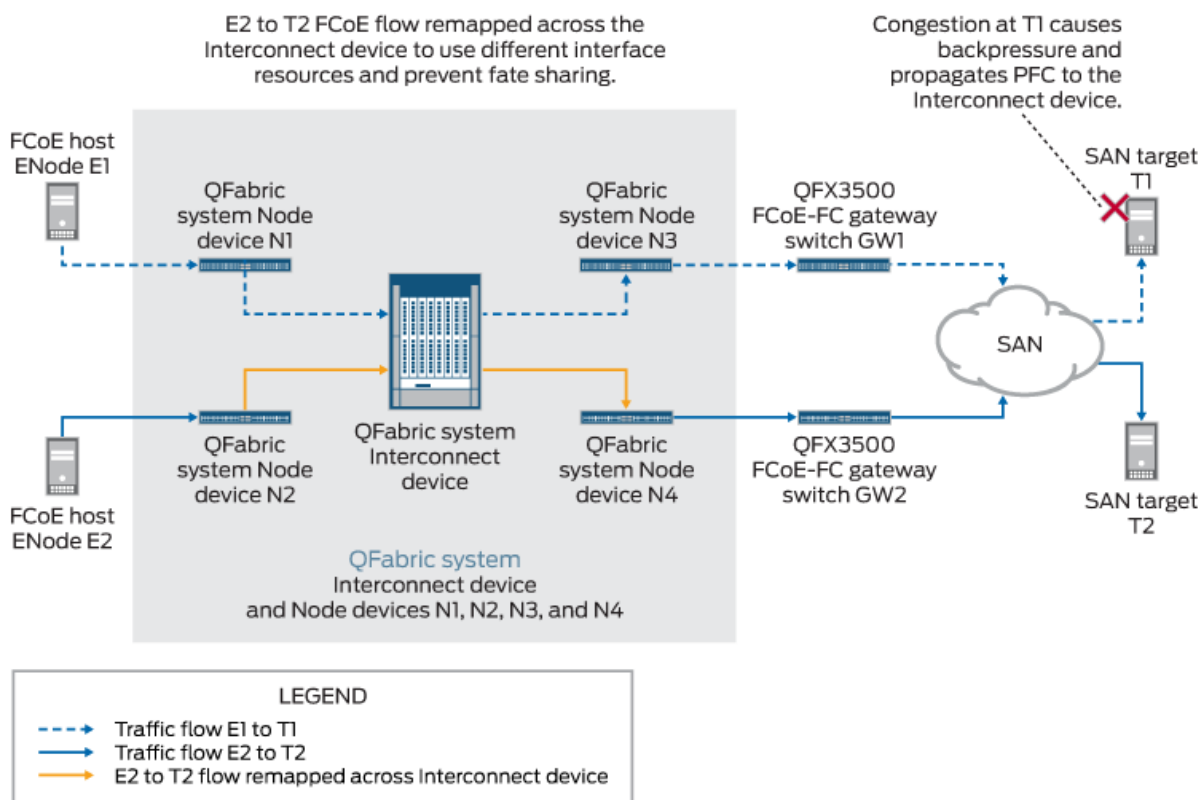


Figure 17 on page 417 is similar to Figure 16 on page 415, with one exception. There are still two FCoE traffic flows, one from host ENode E1 to SAN target T1, and one from host ENode E2 to SAN target T2. Target T1 is experiencing congestion, and target T2 is not experiencing congestion.

The difference is that the path from Node device N2 to the Interconnect device and from the Interconnect device to Node device N4 (yellow in color display, light gray in black and white display) indicates that the forwarding class has been remapped from the original `fcoe` default forwarding class into a different forwarding class.

As in the fate sharing scenario, the flows from hosts E1 and E2 share an ingress interface as they cross the Interconnect device. Also as in the fate sharing scenario, when target T1 experiences congestion, it sends a pause message to temporarily stop the incoming flow until the congestion clears, and the pause message propagates back through the data path.

However, unlike the flows in the fate sharing scenario, these flows use different output queues because the flows have been remapped into different forwarding classes for transit across the Interconnect device. Since the flows use different forwarding classes, they do not share fate on the shared ingress

interface, and the uncongested flow from E2-to-T2 does not share the fate of the congested E1-to-T1 flow.

You configure a firewall filter to control how the forwarding classes are remapped before traffic exits the ingress Node device and crosses the Interconnect device. In the same firewall filter, you also configure terms to control how the remapped forwarding class is mapped back to its original forwarding class when traffic enters the egress Node device after the traffic crosses the Interconnect device. The firewall filter requires terms for remapping forwarding class in both directions of flow. For example, the filters must remap the forwarding class not only in the E1-to-T1 direction, but also in the T1-to-E1 direction.

**NOTE:** If an ENode (FCoE device on the Ethernet network) is directly connected to a QFabric system Node device, and that Node device is directly connected to the FCoE-FC gateway by a LAG interface, then using firewall filters to mitigate fate sharing by remapping forwarding classes is not supported, so that traffic is not remapped.

On Node devices that have directly connected ENodes and that also connect directly to an FCoE-FC gateway using a LAG interface, configure the Node device interfaces in a different VLAN than the interfaces on which you want to mitigate fate sharing. In this scenario, interfaces on the Node device should not be in the same VLAN as interfaces on which you want to apply firewall filters to mitigate fate sharing.

If the interface between the Node device and the FCoE-FC gateway is not a LAG interface, then forwarding class remapping works when ENodes are directly connected to the Node device. The fate sharing mitigation feature does not work only when ENodes are directly connected to the Node device *and* the connection between the Node device and the FCoE-FC gateway is a LAG interface.

## Fate Sharing Mitigation Process

The following sequence summarizes the packet flow and the QFabric system operations for mitigating fate sharing:

1. A packet enters a QFabric system ingress Node device. The ingress Node device classifies the packet into a forwarding class, usually based on its IEEE 802.1p code point (priority).
2. The Node device switching lookup determines that the packet needs to traverse the Interconnect device.
3. On the Node device, a firewall filter remaps the packet from its original forwarding class into a different forwarding class.
4. The packet exits the Node device and enters the Interconnect device, using the new (remapped) forwarding class. On the Interconnect device, the new forwarding class is mapped to a different fabric fc-set, and therefore to a different output queue, than the original forwarding class, so it does

not share fate with traffic in the original forwarding class on ingress interfaces. (Each fabric fc-set maps to a different output queue by default, so placing traffic in a different fabric fc-set allows that traffic to use different output queue bandwidth resources than traffic that is mapped to other fabric fc-sets.) The packet crosses the Interconnect device, and then exits.

5. The packet arrives at the egress Node device. At the Node device ingress interface, the same firewall filter remaps the packet from the new forwarding class back into the original forwarding class.

This process remaps the traffic into a different forwarding class for the journey across the Interconnect device, and then maps the traffic back into its original forwarding class to continue the journey to its destination.

6. The egress Node device forwards the packet toward its destination. Because the packet has been mapped back to its original forwarding class, it once again receives the same CoS treatment as similar traffic that was not remapped across the Interconnect device. It is important to map traffic back to its original forwarding class before forwarding traffic toward its destination because the original forwarding class contains similar traffic, and is configured to support the CoS that the traffic type requires and the destination device expects.

**NOTE:** If you configure non-default forwarding classes and use non-default fabric fc-sets, you must also configure queue scheduling for the new forwarding classes on the Node device and for the non-default fabric fc-sets on the Interconnect device.

Mitigating fate sharing consists of configuration steps that create the necessary forwarding classes and firewall filters, apply the firewall filters to traffic, map the forwarding classes to fabric fc-sets on the Interconnect device, and schedule port bandwidth resources (if needed) for the forwarding classes and fabric fc-sets.

## Forwarding Classes (Node Devices)

If you have only a few flows that you want to separate for transit across the Interconnect device, and the default forwarding classes provide enough separation to avoid fate sharing, you do not need to configure new forwarding classes. There are five default forwarding classes on a QFabric system:

- fcoe—Guaranteed delivery for Fibre Channel over Ethernet (FCoE) traffic.
- no-loss—Guaranteed delivery for TCP lossless traffic.
- best-effort—Provides best-effort delivery without a service profile. Loss priority is typically not carried in a class-of-service (CoS) value.
- network-control—Supports protocol control and is typically high priority.

- `mcast`—Provides service for multdestination (multicast, broadcast, and destination lookup fail) packets.

For example, if you want to separate FCoE traffic into two separate flows on the Interconnect device, and you are not using the `no-loss` forwarding class for other traffic, you can remap some of the FCoE traffic to the `no-loss` forwarding class and leave the rest in the `fcoe` forwarding class. If this provides sufficient separation of flows, you do not need to create new forwarding classes.

Using the existing default forwarding classes has two more time-saving advantages:

1. You do not need to map the forwarding class to a fabric `fc-set` on the Interconnect device, because each default forwarding class is already mapped to a default fabric `fc-set`.
2. You do not need to schedule port bandwidth resources for a new forwarding class on the Node device or for the fabric `fc-set` on the Interconnect device, because each default forwarding class and fabric `fc-set` already has a default port bandwidth allocation.

However, if the default forwarding classes are not sufficient, you can configure up to eight unicast forwarding classes (including the four default forwarding classes) and up to four multdestination forwarding classes (including the default `mcast` forwarding class). You can configure up to six of the unicast forwarding classes as lossless forwarding classes. Lossless transport is not supported on multdestination queues.

For more information about forwarding classes, see ["Understanding CoS Forwarding Classes" on page 124](#). For an example of how to configure forwarding classes, see ["Defining CoS Forwarding Classes" on page 131](#).

**NOTE:** Configuring a new forwarding class includes mapping an output queue to that forwarding class. When you configure a new forwarding class, you also need to configure scheduling resources (output queue bandwidth) for the new forwarding class. When the fate mitigation firewall filter separates traffic flows in a VLAN and assigns all or some of those flows to the new forwarding class, if the output queue mapped to the new forwarding class does not receive bandwidth, traffic cannot be forwarded. For more information about scheduling on Node devices, see ["Understanding CoS Output Queue Schedulers" on page 186](#), ["Understanding CoS Priority Group Scheduling" on page 214](#), and ["Understanding CoS Hierarchical Port Scheduling \(ETS\)" on page 223](#).

## Firewall Filter Construction (Node Devices)

Fate sharing mitigation uses firewall filters to separate traffic before the traffic crosses the Interconnect device, and to bring that traffic back together after it exits the Interconnect device. The QFabric system uses firewall filters to identify (match) traffic and remap forwarding classes because firewall filter match

conditions are granular enough to easily identify and separate (filter) particular traffic flows within a VLAN.

**NOTE:** You can configure firewall filters for fate sharing mitigation only in the firewall family ethernet-switching hierarchy. You cannot configure firewall filters to mitigate fate sharing in the inet (IPv4) or inet6 (IPv6) firewall family hierarchies.

You bind firewall filters for fate sharing mitigation to ingress VLANs as input filters (later in this document is an explanation of why ingress VLANs are the filter bind point). Each firewall filter consists of terms that contain match conditions to identify traffic, and actions to perform on the matched traffic. For more information about firewall filters, see [Overview of Firewall Filters \(QFX Series\)](#).

The firewall filter terms:

1. Remap some or all of the traffic in one forwarding class into another forwarding class before that traffic exits an ingress Node device to go to the Interconnect device. This separates traffic flows before they traverse the Interconnect device, so the traffic uses different output queues and does not experience fate sharing if the traffic in the different forwarding classes uses the same ingress interface.
2. Map the remapped traffic back into its original forwarding class after it exits the Interconnect device, when the traffic enters egress Node device. This brings the traffic flows back into their original forwarding class and classification before the traffic is forwarded toward its destination.

**NOTE:** Each firewall filter requires terms to remap the forwarding class in *both* directions of flow through the Interconnect device. For example, the forwarding class needs to be remapped on the Interconnect device as traffic flows from a source server to a destination target, and the forwarding class also needs to be remapped on the Interconnect device as traffic flows back from the target to the server.

Firewall filter terms contain match conditions (from statement) to identify traffic, and actions (then statement) to tell the system what to do with the identified traffic.

Each forwarding class remapping firewall filter uses match conditions to identify a particular traffic flow to remap, and match conditions to identify the direction of flow on the Interconnect device. Each fate sharing mitigation firewall filter includes terms that:

1. Identify and remap traffic flowing from the server to the target before it enters the Interconnect device.
2. Identify and remap traffic flowing from the server to the target after it exits the Interconnect device.



3. Identify and remap traffic flowing from the target to the server before it enters the Interconnect device.
4. Identify and remap traffic flowing from the target to the server after it exits the Interconnect device.
5. Accept other traffic. Because firewall filters have an implicit default *discard* terminating action, include a final *accept* term so that traffic that does not match the filter is not dropped.

You can use the following match conditions in the filter term from statement to identify traffic that you want to remap as it crosses the Interconnect device:

- Client-side MAC address (for example, an FCF MAC address for FCoE traffic) (`destination-mac-address mac-address`) or (`source-mac-address mac-address`)
- Server-side MAC address (for example, an ENode MAC address for FCoE traffic) (`destination-mac-address mac-address`) or (`source-mac-address mac-address`)
- EtherType (`ether-type value`)

**NOTE:** If you remap an FCoE flow using EtherType as a match condition, you need to include two terms in the filter in each direction of flow to identify the traffic, one term to identify FCoE traffic (EtherType 0x8906), and one term to identify FIP traffic (EtherType 0x8914).

- VLAN (`vlan (vlan-name | vlan-id)`)
- .1q user priority (`dot1q-user-priority value`)

These five match conditions select the traffic from within a VLAN that you want to map to a different forwarding class. The match conditions enable you to identify traffic in VLANs that carry a mix of traffic types—for example, you can identify a flow within a VLAN based on EtherType or .1q value. For more information about match conditions, see [Firewall Filter Match Conditions and Actions \(QFX5100, QFX5110, QFX5120, QFX5200, EX4600, EX4650\)](#).

**BEST PRACTICE:** For FCoE traffic, we recommend that you use the FCF MAC address (instead of the ENode MAC address) as the source or destination address when you configure a firewall filter, because an ENode might be able to reach more than one FCF. Using the FCF MAC is the most specific way to identify the correct path for the traffic.

**NOTE:** You cannot match on multicast addresses based on prefix. You must use a specific multicast address as the source or destination address.

In the same filter term `from` statement, you specify a match condition to determine whether you are identifying traffic that is flowing from a Node device into the Interconnect device, or traffic that is flowing from the Interconnect device to a Node device:

- `to-fabric <except>`—This condition matches traffic that flows from a Node device to an Interconnect device (traffic that is exiting a Node device and entering the Interconnect device). Traffic that matches the `to-fabric` condition is remapped before it exits the ingress Node device and enters the Interconnect device.

The `except` option remaps forwarding classes for traffic that is locally switched. For example, if a target device is directly connected to a Node device, the traffic destined for the directly connected target is remapped to the new forwarding class. When you specify the `except` option, traffic that is remotely switched is *not* remapped to a new forwarding class before it crosses the Interconnect device.

- `from-fabric`—This condition matches traffic that flows from the Interconnect device to a Node device (traffic that is exiting the Interconnect device and entering the egress Node device). Traffic that matches the `from-fabric` condition is mapped back to its original forwarding class after it exits the Interconnect device, when it enters the egress Node device.

**BEST PRACTICE:** In a firewall filter configuration, if you use a `to-fabric except` match condition, place it before the `from-fabric` term in the sequence of terms in the filter.

After you configure match conditions in a filter term, you configure an action to take on the identified (matched) traffic in the same term. Because the goal is to remap traffic in one forwarding class into a different forwarding class, the action is usually to place the matched traffic into a forwarding class.

Use the following actions (`then` statement) to control into which forwarding class the matched traffic is remapped in a given term:

- `forwarding-class forwarding-class-name`—Specify a default or a user-defined forwarding class into which matching traffic is mapped.
- `loss-priority level`—If you specify a forwarding class for matching traffic, you must also specify the packet loss priority (PLP) level for the forwarding class. The PLP level can be `low`, `medium-high`, or `high`.
- `count counter-name`—Optionally, you can configure an action to count the number of packets affected by each term.

**NOTE:** You can use the match conditions to identify a traffic flow, and then count the packets without remapping the forwarding class. To do that, in the `then` statement, do not include the forwarding class and loss priority, include only the `count` action.

## Applying Firewall Filters to Traffic (Node Devices)

You apply (bind) firewall filters for fate sharing mitigation to ingress VLANs, not to ports. (Firewall filters for mitigating fate sharing do not apply to VLANs on the egress side.) Applying the firewall filter to an ingress VLAN has advantages compared to applying the firewall filter to a port:

- The filter affects all of the matched traffic on all interfaces that are members of the VLAN, on all Node devices on the QFabric system. Instead of applying the firewall filter to individual ports or ranges of ports on each Node device, you only have to apply the firewall filter once to the VLAN.
- VLANs usually carry similar types of traffic.

You bind firewall filters to ingress VLANs as input filters using the `set vlans vlan-name filter input filter-name` configuration statement. See [Configuring Firewall Filters](#) for more information about configuring and applying firewall filters.

**BEST PRACTICE:** Place traffic of one type in one VLAN (use separate VLANs for each different type of traffic). We recommend that you do not mix different types of traffic in the same VLAN. The QFabric system requires that a VLAN that carries FCoE traffic must carry only FCoE traffic. However, it is a good practice to do the same thing with other types of traffic. For example, if your network carries both iSCSI and NAS traffic, we recommend that you dedicate one VLAN to iSCSI traffic, and one VLAN to NAS traffic (and so on). You can configure separate firewall filters to mitigate fate sharing for each type of traffic.

**NOTE:** Because firewall filters for mitigating fate sharing are applied to VLANs, and not to ports, there are several behaviors you should be aware of:

- If more than one VLAN uses a port, the firewall filter applies only to the traffic in the VLAN on which you applied the firewall filter. Traffic in other VLANs might be exposed to fate sharing on the Interconnect device.
- The ports on which the firewall filter is applied depend on VLAN membership. If ports on multiple Node devices are members of the VLAN, then the firewall filter remaps traffic on the VLAN member ports of all of those Node devices. If you want to remap traffic on only one Node device, then the VLAN member interfaces should all be on that Node device, and not on other Node devices. (Configuring a VLAN that includes member interfaces from only one Node device enables you to remap traffic on that Node device independently from other Node devices.)

- Although firewall filters mitigate fate sharing on the Interconnect device, they do not mitigate fate sharing on a Node device. This is because PFC is applied to specified queues on a port, not to a VLAN. (Recall that forwarding classes are mapped to queues, so all traffic in the same forwarding class uses the same queue, regardless of VLAN membership.)

An example scenario is two VLANs that contain FCoE traffic that is classified into the `fcoe` forwarding class and use an ingress interface on the same Node device. The `fcoe` forwarding class is classified to IEEE 802.1p code point 011 (priority 3) to identify the FCoE traffic on both VLANs (because all of the FCoE traffic requires the same CoS treatment and all of the traffic is destined for the same SAN), and so both VLANs use the same output queue.

If FCoE traffic in one of the VLANs experiences congestion, PFC is enabled on the flow, and the flow is paused until the congestion clears. Because the FCoE traffic in the other VLAN uses the same output queue (forwarding class), when the congested FCoE flow is paused on the ingress interface, all FCoE traffic that uses that ingress interface is also paused. In this way, the congested FCoE flow affects the uncongested FCoE flow, and the two flows share the same fate.

So if two servers on the same Node device ingress port send traffic that belongs to the same forwarding class (in this example, `fcoe`), they can experience fate sharing on the Node device.



**WARNING:** Do not apply firewall filters that remap forwarding classes while traffic that the filters affect is flowing!

For forwarding class remapping to work properly, traffic must be mapped from its original forwarding class to a new forwarding class before it enters the Interconnect device, and then mapped back to the original forwarding class after it exits the Interconnect. If traffic is not mapped back into its original forwarding class after crossing the Interconnect device, traffic is classified into the wrong forwarding class and is not delivered as expected. Because of this, the QFabric system must program the filters on the ingress Node device and the egress Node device when affected traffic is not flowing.

If traffic is flowing when you apply the filters to a VLAN, and the ingress Node device filter is programmed before the egress Node device filter is programmed, traffic is not remapped back into its original forwarding class until the egress Node device filter is applied. For this reason, apply filters only when affected traffic is not flowing through the QFabric system.

## Mapping Forwarding Classes to Fabric Forwarding Class Sets (Interconnect Device)

The five default forwarding classes (best-effort, fcoe, no-loss, network-control, and mcast, see ["Forwarding Classes \(Node Devices\)" on page 419](#)) are mapped by default to fabric fc-sets on the Interconnect device. If you are using only default forwarding classes on the Node devices, then you do not need to map forwarding classes to fabric fc-sets, you can use the default mapping.

If you create new (user-defined) forwarding classes on a Node device, you must map the new forwarding classes to fabric fc-sets on the Interconnect device. (If you do not map a new forwarding class to a fabric fc-set on the Interconnect device, the traffic that belongs to the new forwarding class receives very little bandwidth on the Interconnect device.)

Each fabric fc-set maps to a different output queue on the Interconnect device by default, much like each forwarding class maps to a different output queue by default on Node devices. Mapping a new forwarding class to a non-default (unused) fabric fc-set causes the traffic assigned to that forwarding class to use a different output queue on the Interconnect device. (The traffic in the new forwarding class uses a different output queue than traffic mapped to other fabric fc-sets.)

Also similar to forwarding classes on Node devices, there are five default fabric fc-sets on the Interconnect device, and twelve total fabric fc-sets, eight of which are unicast fabric fc-sets, and four of which are multidestination fabric fc-sets. Each default forwarding class has a default mapping to one of the default fabric fc-sets. The non-default fabric fc-sets are hidden until you map forwarding classes to them, but are available for use.

The five default forwarding classes are mapped to the five default fabric fc-sets as shown in [Table 83 on page 426](#) (you can reconfigure the mapping of default forwarding classes to default fabric fc-sets if you want):

**Table 83: Default Fabric Forwarding Class Sets**

| Fabric Forwarding Class Set Name | Characteristics                                                                                                                                                                                                                                                                                                                                                                                                                                    |
|----------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| fabric_fcset_be                  | Transports best-effort unicast traffic across the fabric.                                                                                                                                                                                                                                                                                                                                                                                          |
| fabric_fcset_strict_high         | <p>Transports unicast traffic that has been configured with strict-high priority and in the network-control forwarding class across the fabric.</p> <p><b>NOTE:</b> This fabric fc-set receives as much bandwidth across the fabric as it needs to service the traffic in the group up to the entire fabric interface bandwidth. For this reason, exercise caution when mapping traffic to this fabric fc-set to avoid starving other traffic.</p> |

**Table 83: Default Fabric Forwarding Class Sets (Continued)**

| Fabric Forwarding Class Set Name     | Characteristics                                                                                                                                              |
|--------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------|
| <code>fabric_fcset_noloss1</code>    | Transports unicast traffic in the default fcoe forwarding class across the fabric.                                                                           |
| <code>fabric_fcset_noloss2</code>    | Transports unicast traffic in the default no-loss forwarding class across the fabric.                                                                        |
| <code>fabric_fcset_multicast1</code> | Transports multdestination traffic in the mcast forwarding class across the fabric. This fabric fc-set is valid only for multdestination forwarding classes. |

The remaining four unicast fabric fc-sets (`fabric_fcset_noloss3`, `fabric_fcset_noloss4`, `fabric_fcset_noloss5`, and `fabric_fcset_noloss6`) can carry lossless traffic and are available for mapping or remapping forwarding classes on the Interconnect device. The remaining three multdestination fabric fc-sets (`fabric_fcset_multicast2`, `fabric_fcset_multicast3`, and `fabric_fcset_multicast4`) are available for remapping multdestination forwarding classes.

The total of six lossless and four multdestination fabric fc-sets enable you to separate traffic from Node devices into up to ten classes on the Interconnect device, not including the best-effort and strict high-priority fabric fc-sets. Because each fabric fc-set uses a different output queue on egress interfaces, the flows (forwarding classes) mapped to each fabric fc-set do not share fate with flows in other fabric fc-sets on ingress interfaces.

The total number of unique flows on a QFabric system is vastly greater than the number of fabric fc-sets, so fabric fc-sets still aggregate flows—each fabric fc-set will carry a group of flows that require similar CoS treatment. However, the fabric fc-sets enable you to spread the flows across multiple output queues, and thus mitigate the effects of fate sharing.

**NOTE:** The forwarding class flows within a fabric fc-set share fate on ingress interfaces because they use the same output queue. However, the ability to separate flows into different classes that use different output queues enables you to control how much traffic is mapped to a given output queue, and in that way to mitigate the possibility of fate sharing.

For more information about fabric fc-sets, see ["Understanding CoS Fabric Forwarding Class Sets" on page 396](#).

## Scheduling Bandwidth for Fabric Forwarding-Class Sets (Interconnect Device)

The five default fabric fc-sets (`fabric_fcset_be`, `fabric_fcset_strict_high`, `fabric_fcset_noloss1`, `fabric_fcset_noloss2`, and `fabric_fcset_multicast1`, see ["Mapping Forwarding Classes to Fabric Forwarding Class Sets \(Interconnect Device\)" on page 426](#)) receive scheduling resources on Interconnect device output queues by default. If you are using only default fabric fc-sets, then you can use the default scheduling. However, you can change scheduling parameters, such as the amount of bandwidth allocated to a default fabric fc-set, if you want to adjust the default scheduling.

If you configure a new forwarding class on a Node device, you must map the new forwarding class to a fabric fc-set so that the traffic classified into the forwarding class receives queue bandwidth resources. If you map a new forwarding class to one of the default fabric fc-sets on the Interconnect device, then the default bandwidth scheduled for that fabric fc-set is shared among the forwarding classes assigned to the fabric fc-set by default, and also with the new forwarding class.

If you map a new forwarding class to one of the non-default fabric fc-sets, you must schedule queue bandwidth resources for that fabric fc-set, or else the traffic mapped to the fabric fc-set receives only a small amount of bandwidth.

**NOTE:** You apply queue (forwarding class) scheduling to interfaces. The Interconnect device interfaces consist of the ingress and egress 40-Gbps (*fte*) interfaces that connect to QFabric system Node devices, and internal Clos fabric (*bfte*) interfaces. You need to apply the appropriate scheduler to each *fte* interface in the traffic path. All traffic traverses the internal Clos fabric interfaces, so you also need to apply the appropriate scheduler to the Clos fabric *bfte* interfaces. (You configure one scheduler that applies to all of the internal Clos fabric interfaces. It is not possible or desirable to attach a scheduler to a particular internal Clos fabric interface.) Because one scheduler applies to all of the Clos fabric interfaces, you either use the default scheduler on all Clos interfaces, or you use your custom configured scheduler on all Clos interfaces.

For conceptual information about configuring CoS scheduling on an Interconnect device and across the entire QFabric system, see ["Understanding CoS Scheduling Across the QFabric System" on page 327](#). For information about default CoS scheduling on the Interconnect device, see ["Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\)" on page 315](#). For an example of how to configure scheduling on an Interconnect device and across the entire QFabric system, see ["Example: Configuring CoS Scheduling Across the QFabric System" on page 353](#).

## Multidestination Traffic (FCoE Initialization Protocol Traffic)

Multidestination (multicast, broadcast, and destination lookup fail) traffic that is not switched locally on a QFabric system Node device is switched across the Interconnect device. On the Node device, by

default, multdestination traffic uses the `mcast` forwarding class. On the Interconnect device, by default the multdestination traffic from the Node devices uses the `fabric_fcset_multicast1` fabric fc-set. The output queue for the `fabric_fcset_multicast1` fabric fc-set receives up to 20 percent of the available egress port bandwidth.

FCoE devices on the Ethernet network use FCoE Initialization Protocol (FIP) to establish a virtual point-to-point link with the FCF. The FCF sends periodic multicast discovery advertisements (MDAs) to advertise its presence on the network to ENodes. When an ENode comes online, it sends a multicast discovery solicitation (MDS) advertisement to search the network for FCFs.

The FIP MDA and MDS advertisements use the default multicast queue on the Interconnect device. If the amount of multdestination traffic that crosses the Interconnect device causes congestion on the multdestination queue, that congestion can impact FIP discovery advertisement traffic. (Fate sharing can occur because the FIP advertisements share the same fabric fc-set, and therefore the same output queue, as the rest of the multdestination traffic on the Interconnect device. Multdestination traffic that uses the same ingress interface at any point on the Interconnect device can experience fate sharing if the output queue becomes congested.)

**NOTE:** If the amount of multdestination traffic on the Interconnect device is not enough to cause congestion, you do not have to remap multicast FIP traffic into a separate forwarding class to avoid fate sharing.

**NOTE:** Although multicast FIP traffic uses the `mcast` queue and the `fabric_fcset_multicast1` fabric fc-set by default, unicast FCoE and FIP traffic uses the `fcoe` forwarding class and the `fabric_fcset_noloss1` fabric fc-set by default.

If the amount of multdestination traffic that traverses the Interconnect device can cause congestion, then you can remap the FIP multicast traffic into a new forwarding class on the Node device and a new fabric fc-set on the Interconnect device to mitigate fate sharing. The process is similar to mitigating fate sharing on unicast traffic, but there are a few differences:

1. Configure a new multdestination forwarding class for the FIP multicast traffic on the Node device. (By default, multicast FIP traffic is classified into the default `mcast` forwarding class.)
2. Configure queue and priority group scheduling (hierarchical scheduling) for the new multdestination forwarding class.
3. Configure a firewall filter to remap the FIP multicast traffic into the new forwarding class. To match FIP multicast traffic, specify two match conditions: the ALL-FCF-MAC address (01:10:18:01:00:02) as the source or destination MAC address (depending on the direction of flow), and the FIP EtherType (0x8014).



4. Bind the firewall filter to the appropriate VLAN.
5. Map the new multidestination forwarding class that you created on the Node device to an unused multicast fabric fc-set on the Interconnect device.
6. Configure scheduling for the multicast fabric fc-set on the Interconnect device.

**NOTE:** When configuring firewall filter match conditions, you cannot match on multicast addresses based on prefix. You must use a specific multicast address as the source or destination address.

## Best Practices

The previous sections include some best practices for mitigating fate sharing. This section aggregates those best practices, along with a few other tips.

### VLANs

Place traffic of one type in one VLAN (use separate VLANs for each different type of traffic). We recommend that you do not mix different types of traffic in the same VLAN. The QFabric system requires that a VLAN that carries FCoE traffic must carry only FCoE traffic. However, it is a good practice to do the same thing with other types of traffic. For example, if your network carries both iSCSI and NAS traffic, we recommend that you dedicate one VLAN to iSCSI traffic, and one VLAN to NAS traffic (and so on). You can configure separate firewall filters to mitigate fate sharing for each type of traffic.

### Source/Destination MAC Address for FCoE Traffic

For FCoE traffic, we recommend that you use the FCF MAC address (instead of the ENode MAC address) as the source or destination address when you configure a firewall filter, because an ENode might be able to reach more than one FCF. Using the FCF MAC is the most specific way to identify the correct path for the traffic.

### Firewall Filter Term Sequence

In most cases, the sequence of terms in a fate sharing firewall filter does not matter (with the exception of the final accept term), so in most cases, it does not matter if a `from-fabric` term is placed before a `to-fabric` term in the firewall filter.

However, we recommend that if you use the `except` option with `to-fabric` (`to-fabric except`), you should place the `to-fabric except` term before the `from-fabric` in the firewall filter.

In general, we recommend that in a filter, you configure the `to-fabric` terms first, then configure the `from-fabric` terms, and end the filter with an `accept` term (unless you want to drop traffic that does not match the filter).

## Limitations and Notes on Behavior

There are a number of limitations and behaviors that you should understand about how to mitigate fate sharing across an Interconnect device. Some of those limitations and behaviors have been discussed in the previous sections, and are repeated here for your convenience.

### Limitations

- You can configure firewall filters for fate sharing mitigation only in the `firewall family ethernet-switching` hierarchy. You cannot configure firewall filters to mitigate fate sharing in the `inet (IPv4)` or `inet6 (IPv6)` firewall family hierarchies.
- Interconnect device fabric `fc-sets` are not user-configurable (you cannot rename them or configure new fabric `fc-sets`). You can map the default forwarding classes and the forwarding classes that you define on Node devices to fabric `fc-sets` to control the traffic that is mapped to each fabric-`fcset`.
- The possibility of fate sharing cannot be avoided on Node device ingress interfaces. If two servers access a Node device on the same ingress interface, and both servers send traffic flows that are classified into the same forwarding class (for example, if both flows are FCoE traffic and are classified into the `fcqe` forwarding class), then even if the flows are in different VLANs, congestion on one flow affects the other flow. The congested flow affects the uncongested flow because both flows share the same forwarding class (and therefore the same output queue and IEEE 802.1p code point), and priority-based flow control (PFC) is applied to the ingress interface, not to the VLAN. So when PFC pauses the congested flow on the ingress interface, the uncongested flow that uses the same code point is also paused.
- The Interconnect device supports a maximum of six lossless unicast flow groups (six lossless unicast fabric `fc-sets`). In practice, a QFabric system has many more than six flows, so you cannot map each individual flow to a dedicated fabric `fc-set`. However, you can group flows into six separate sets by mapping groups of flows to different fabric `fc-sets`. Each fabric `fc-set` uses a different output queue, so the flows in one fabric `fc-set` do not share fate with the flows in the other fabric `fc-sets` when the flows traverse the same ingress interface. The ability to separate flows into six different fabric `fc-sets` spreads the flows among six different output queues, thus mitigating fate sharing.
- The Interconnect device supports a maximum of four multidestination flow groups (four multicast fabric `fc-sets`).
- The flows (forwarding classes) within a fabric `fc-set` share fate when they use the same ingress interface because they use the same output queue. (However, the ability to separate flows into

different classes that use different output queues enables you to control how much traffic is mapped to a given output queue, and to mitigate the possibility of fate sharing.)

- Do not apply firewall filters that remap forwarding classes while traffic that the filters affect is flowing!

For forwarding class remapping to work properly, traffic must be mapped from its original forwarding class to a new forwarding class before it enters the Interconnect device, and then mapped back to the original forwarding class after it exits the Interconnect. If traffic is not mapped back into its original forwarding class after crossing the Interconnect device, traffic is classified into the wrong forwarding class and is not delivered as expected. Because of this, the QFabric system must program the filters on the ingress Node device and the egress Node device when affected traffic is not flowing.

If traffic is flowing when you apply the filters to a VLAN, and the ingress Node device filter is programmed before the egress Node device filter is programmed, traffic is not remapped back into its original forwarding class until the egress Node device filter is applied. For this reason, apply filters only when affected traffic is not flowing through the QFabric system.

- If an ENode (FCoE device on the Ethernet network) is directly connected to a QFabric system Node device, and that Node device is directly connected to the FCoE-FC gateway by a LAG interface, then using firewall filters to mitigate fate sharing by remapping forwarding classes is not supported, so that traffic is not remapped.

On Node devices that have directly connected ENodes and that also connect directly to an FCoE-FC gateway using a LAG interface, configure the Node device interfaces in a different VLAN than the interfaces on which you want to mitigate fate sharing. In this scenario, interfaces on the Node device should not be in the same VLAN as interfaces on which you want to apply firewall filters to mitigate fate sharing.

If the interface between the Node device and the FCoE-FC gateway is not a LAG interface, then forwarding class remapping works when ENodes are directly connected to the Node device. The fate sharing mitigation feature does not work only when ENodes are directly connected to the Node device *and* the connection between the Node device and the FCoE-FC gateway is a LAG interface.

- When configuring firewall filter match conditions, you cannot match on multicast addresses based on prefix. You must use a specific multicast address as the source or destination address.

## Notes on Behavior

- You bind (apply) firewall filters for mitigating fate sharing to ingress VLANs only, not to ports. The filter affects all matched traffic on all Node device ingress interfaces that are members of the VLAN. So if ports on multiple Node devices are members of the VLAN, then the firewall filter remaps traffic on the VLAN member ports of all of those Node devices. If you want to remap traffic on only one

Node device, then the VLAN member interfaces should all be on that Node device, and not on other Node devices.

- Although firewall filters mitigate fate sharing on the Interconnect device, they do not mitigate fate sharing on a Node device. This is because PFC is applied to specified queues on a port, not to a VLAN. (Recall that forwarding classes are mapped to queues, so all traffic in the same forwarding class uses the same queue, regardless of VLAN membership.)

An example scenario is two VLANs that contain FCoE traffic that is classified into the `fcoe` forwarding class and use an ingress interface on the same Node device. The `fcoe` forwarding class is classified to IEEE 802.1p code point 011 (priority 3) to identify the FCoE traffic on both VLANs (because all of the FCoE traffic requires the same CoS treatment and all of the traffic is destined for the same SAN), and so both VLANs use the same output queue.

If FCoE traffic in one of the VLANs experiences congestion, PFC is enabled on the flow, and the flow is paused until the congestion clears. Because the FCoE traffic in the other VLAN uses the same output queue (forwarding class), when the congested FCoE flow is paused on the ingress interface, all FCoE traffic that uses that ingress interface is also paused. In this way, the congested FCoE flow affects the uncongested FCoE flow, and the two flows share the same fate.

So if two servers on the same Node device ingress port send traffic that belongs to the same forwarding class (in this example, `fcoe`), they can experience fate sharing on the Node device.

- When you configure a firewall filter, by default, the last term in the filter is a discard action. (This is standard default behavior and is not unique to fate sharing mitigation filters.) To avoid dropping traffic that does not match the filter conditions for forwarding class remapping, add a final term with `accept` as the action. This is especially important when you are not remapping all of the traffic in a VLAN.
- If you remap FCoE flows based on EtherType, include separate filter terms to match both the FCoE EtherType (0x8906) and the FIP EtherType (0x8914).
- You must configure filter terms that remap the forwarding classes in both directions of flow. You need to configure terms for `to-fabric` and `from-fabric` for the flow from the originating device to the target, and also for the return flow from the target to the originating device. For example, for an FCoE flow, you configure a `to-fabric` and a `from-fabric` term for the traffic flowing from the ENode to the FC SAN, and a `to-fabric` and a `from-fabric` term for traffic flowing from the FC SAN to the ENode.

## RELATED DOCUMENTATION

*Overview of Firewall Filters (QFX Series)*

*Firewall Filter Match Conditions and Actions (QFX5100, QFX5110, QFX5120, QFX5200, EX4600, EX4650)*

*Configuring Firewall Filters*

|                                                                                                                        |     |
|------------------------------------------------------------------------------------------------------------------------|-----|
| Understanding CoS Fabric Forwarding Class Sets                                                                         | 396 |
| Understanding CoS Scheduling Across the QFabric System                                                                 | 327 |
| Understanding Default CoS Scheduling on QFabric System Interconnect Devices (Junos OS Release 13.1 and Later Releases) | 315 |
| Example: Configuring CoS Scheduling Across the QFabric System                                                          | 353 |
| Understanding CoS Output Queue Schedulers                                                                              | 186 |
| Understanding CoS Hierarchical Port Scheduling (ETS)                                                                   | 223 |
| Understanding CoS Forwarding Classes                                                                                   | 124 |
| Defining CoS Forwarding Classes                                                                                        | 131 |
| Configuring Fate Sharing Mitigation Across the Interconnect Device by Remapping Traffic Flows (Forwarding Classes)     | 434 |

## Configuring Fate Sharing Mitigation Across the Interconnect Device by Remapping Traffic Flows (Forwarding Classes)

On a QFabric system, traffic flows that belong to the same forwarding class are mapped to the same output queue and share the output queue resources. If congestion occurs on one of these flows, the congestion can affect the uncongested flows in the forwarding class when the flows use the same ingress interface.

For example, if a congested flow is paused to prevent packet loss, uncongested flows that use the same ingress interface are also paused because they share the same forwarding class and output queue. When a congested flow affects an uncongested flow, the flows share the same fate—this is known as fate sharing.

Fate sharing happens because pausing traffic is based on forwarding class. When a flow experiences congestion, the output queue sends a pause message to the input queue on which the flow arrived. On that input queue, the pause message affects all traffic in the forwarding class that is mapped to the congested output queue. So all traffic in that forwarding class is paused on the input queue, not just the flow that is experiencing the congestion. This is how uncongested flows can share fate with a congested flow.

Traffic from many QFabric system Node devices crosses the Interconnect device, so flows within a given forwarding class are aggregated on the Interconnect device. The aggregated flows use the same output queue on the Interconnect device and are subject to fate sharing if the flows also use the same ingress interface.

In addition to the external physical interfaces that connect the Interconnect device to Node devices, the Interconnect device has internal Clos interfaces. The Interconnect device automatically selects the best path through its internal Clos interfaces. Path selection through the internal Clos interfaces is not

configurable, so you cannot control the traffic that enters any particular ingress Clos interface, and so fate sharing can occur on the Interconnect device. (On Node devices, you control the traffic connected to an ingress interface, but on the Interconnect device, you cannot control which flows use a particular internal ingress Clos interface.)

However, you can use firewall filters to separate the traffic assigned to one forwarding class and split it into different forwarding classes for the journey across the Interconnect device. Remapping the flows into different forwarding classes means the flows use different output queues on the Interconnect device. If the flows use the same ingress interface on the Interconnect device, they do not experience fate sharing because only the flows mapped to the congested queue are paused, while the flows remapped to other forwarding classes are not paused.

This topic shows you how to configure firewall filters to remap traffic across the Interconnect device and mitigate fate sharing.

To change the forwarding class (and therefore the output queue) that traffic uses on the Interconnect device, you need to map traffic into a new forwarding class before it enters the Interconnect device, then map the traffic back into the original forwarding class after it exits the Interconnect device. Traffic needs to be mapped back into its original forwarding class before it leaves the QFabric system because the original forwarding class contains similar traffic, and is configured to support the CoS that the traffic type requires and the destination device expects. For example, FCoE traffic destined for different targets in the same Fibre Channel storage area network must be in the same forwarding class (and therefore have the same IEEE 802.1p priority), or the traffic is not handled properly.

The firewall filter has to remap traffic in both directions of flow. For example, if a flow transports traffic between a server and a target device, remapping needs to occur when traffic flows from the server to the target device, and also when traffic flows from the target device to the server. Firewall filter terms contain match conditions (*from* statement) to identify traffic, and actions (*then* statement) to tell the system what to do with the identified traffic.

You configure a firewall filter for fate sharing mitigation in the `firewall family ethernet-switching` hierarchy. You cannot configure firewall filters to mitigate fate sharing in the `inet` (IPv4) or `inet6` (IPv6) firewall family hierarchies.

To mitigate fate sharing across the Interconnect device, you need to configure a firewall filter that:

1. Identifies and remaps traffic flowing from a source to a destination before it enters the Interconnect device. (This separates flows for crossing the Interconnect device.)
2. Identifies and remaps traffic flowing from a source to a destination after it exits the Interconnect device. (This brings flows back into their original forwarding class before traffic is forwarded toward its destination.)

Steps 1 and 2 combine to remap flows across the Interconnect device as traffic travels from a source to a destination.

3. Identifies and remaps traffic flowing back from a destination to a source before it enters the Interconnect device. (This separates flows for crossing the Interconnect device in the other direction.)
4. Identifies and remaps traffic flowing back from a destination to a source after it exits the Interconnect device. (This brings flows back into their original forwarding class in the other direction.)

Steps 3 and 4 combine to remap flows across the Interconnect device on the return path, as traffic flows from the destination device back to the original source device.

5. Accept other traffic. Because firewall filters have an implicit default *discard* terminating action, include a final *accept* term so that traffic that does not match the filter is not dropped (unless you want to drop traffic that does not match the filter).

You can use the following match conditions in the filter term `from` statement to identify (select) traffic that you want to remap as it crosses the Interconnect device:

- Client-side MAC address (for example, an FCF MAC address for FCoE traffic) (`destination-mac-address mac-address`) or (`source-mac-address mac-address`)
- Server-side MAC address (for example, an ENode MAC address for FCoE traffic) (`destination-mac-address mac-address`) or (`source-mac-address mac-address`)
- EtherType (`ether-type value`)

**NOTE:** If you remap an FCoE flow using EtherType as a match condition, you need to include two terms in the filter in each direction of flow to identify the traffic, one term to identify FCoE traffic (EtherType 0x8906), and one term to identify FIP traffic (EtherType 0x8914).

- VLAN (`vlan (vlan-name | vlan-id)`)
- .1q user priority (`dot1q-user-priority value`)

Match conditions enable you to identify traffic in VLANs that carry a mix of traffic types—for example, you can identify a flow within a VLAN based on EtherType or .1q value. For more information about match conditions, see [Firewall Filter Match Conditions and Actions \(QFX5100, QFX5110, QFX5120, QFX5200, EX4600, EX4650\)](#).

**BEST PRACTICE:** For FCoE traffic, we recommend that you use the FCF MAC address (instead of the ENode MAC address) as the source or destination address when you configure a firewall filter, because an ENode might be able to reach more than one FCF. Using the FCF MAC is the most specific way to identify the correct path for the traffic.

**NOTE:** You cannot match on multicast addresses based on prefix. You must use a specific multicast address as the source or destination address.

In the same filter term `from` statement, you specify a match condition to determine whether you are identifying traffic that is flowing from a Node device into the Interconnect device, or traffic that is flowing from the Interconnect device to a Node device:

- `to-fabric <except>`—This condition matches traffic that flows from a Node device to an Interconnect device (traffic that is exiting a Node device and entering the Interconnect device). Traffic that matches the `to-fabric` condition is remapped before it exits the ingress Node device and enters the Interconnect device.

The `except` option remaps forwarding classes for traffic that is locally switched. For example, if a target device is directly connected to a Node device, the traffic destined for the directly connected target is remapped to the new forwarding class. When you specify the `except` option, traffic that is remotely switched is *not* remapped to a new forwarding class before it crosses the Interconnect device.

- `from-fabric`—This condition matches traffic that flows from the Interconnect device to a Node device (traffic that is exiting the Interconnect device and entering the egress Node device). Traffic that matches the `from-fabric` condition is mapped back to its original forwarding class after it exits the Interconnect device, when it enters the egress Node device.

**BEST PRACTICE:** In a firewall filter configuration, if you use a `to-fabric except` match condition, place it before the `from-fabric` term in the sequence of terms in the filter. In general, we recommend that in a filter, you configure the `to-fabric` terms first, then configure the `from-fabric` terms.

After you configure match conditions in a filter term, you configure an action to take on the identified (matched) traffic in the same term. Because the goal is to remap traffic in one forwarding class into a different forwarding class, the action is usually to place the matched traffic into a forwarding class.

Use the following actions (`then` statement) to control into which forwarding class the matched traffic is remapped in a given term:

- `forwarding-class forwarding-class-name`—Specify a default or a user-defined forwarding class into which matching traffic is mapped.
- `loss-priority level`—If you specify a forwarding class for matching traffic, you must also specify the packet loss priority (PLP) level for the forwarding class. The PLP level can be `low`, `medium-high`, or `high`.



- `count counter-name`—Optionally, you can configure an action to count the number of packets affected by each term.

**NOTE:** You can use the match conditions to identify a traffic flow, and then count the packets without remapping the forwarding class. To do that, in the `then` statement, do not include the forwarding class and loss priority, include only the `count` action.

After you configure a firewall filter that remaps traffic across the Interconnect device in both directions of flow, you bind (apply) the filter to an ingress (input) VLAN. The filter only affects traffic in that VLAN.

The following procedure shows how to configure a firewall filter that mitigates fate sharing on the Interconnect device using the CLI. Steps 1-4 configure forwarding class remapping for traffic leaving an ingress Node device and entering the Interconnect device (`to-fabric`), in both directions of flow. Steps 5-8 configure forwarding class remapping for traffic leaving the Interconnect device and entering the egress Node device (`from-fabric`), in both directions of flow.

1. Name the firewall filter and the first term of the filter, and then define match conditions for traffic flowing from the ingress Node device to the Interconnect device in the server-to-target direction (this filter term identifies the traffic to map into a different forwarding class):

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name from flow-match-conditions
user@switch# set filter filter-name term term-name to-fabric
```

The *flow-match-conditions* specify the traffic that you want to remap to a different forwarding class on the ingress Node device for transport across the Interconnect device. The `to-fabric` condition matches only traffic that is going from the ingress Node device to the Interconnect device.

2. In the same firewall filter and term, configure the action to take on traffic on the ingress Node device that matches the conditions in the server-to-target direction (the action is to map the traffic into a different forwarding class on the ingress Node device, before the traffic enters the Interconnect device):

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name then forwarding-class new-forwarding-class-name loss-priority priority-value
user@switch# set filter filter-name term term-name then counter counter-name
```

The *new-forwarding-class-name* specifies the forwarding class that the matching traffic is mapped to for transport across the Interconnect device. The packet counter action is optional, but is included

here and in later steps because many administrators like to have this type of information available to analyze traffic patterns.

3. In the same firewall filter, configure a second term to define match conditions for traffic flowing from the ingress Node device to the Interconnect device in the target-to-server direction (this filter term identifies the traffic to map into a different forwarding class):

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name from flow-match-conditions
user@switch# set filter filter-name term term-name to-fabric
```

The *flow-match-conditions* specify the traffic that you want to remap to a different forwarding class on the ingress Node device for transport across the Interconnect device. The to-fabric condition matches only traffic that is going from the ingress Node device to the Interconnect device.

4. In the second term in the same firewall filter, configure the action to take on traffic on the ingress Node device that matches the conditions in the target-to-server direction (the action is to map the traffic into a different forwarding class on the ingress Node device, before the traffic enters the Interconnect device):

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name then forwarding-class new-forwarding-class-name loss-priority priority-value
user@switch# set filter filter-name term term-name then counter counter-name
```

The first four steps of this process configure match conditions to identify Interconnect device ingress traffic in both directions of flow, and the forwarding class remapping action to take on the matched traffic. The next four steps map the traffic back into its original forwarding class after the traffic exits the Interconnect device, in both directions of flow.

5. In the same firewall filter, configure a third term to define match conditions for traffic flowing from the Interconnect device to the egress Node device in the server-to-target direction (this term identifies traffic to map back into the original forwarding class after it crosses the Interconnect device):

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name from flow-match-conditions
user@switch# set filter filter-name term term-name from-fabric
```

The *flow-match-conditions* specify the traffic that you want to map back into the original forwarding class on the egress Node device, after the traffic crosses the Interconnect device. The from-fabric

condition matches only traffic that is coming from the Interconnect device into the egress Node device.

6. In the third term in the same firewall filter, configure the action to take on traffic when it enters the egress Node device from the Interconnect device in the server-to-target direction (the action is to map the traffic back into its original forwarding class on the egress Node device, after the traffic crosses the Interconnect device):

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name then forwarding-class original-forwarding-class-name loss-priority priority-value
user@switch# set filter filter-name term term-name then counter counter-name
```

The *original-forwarding-class-name* specifies the original forwarding class name (the forwarding class the traffic was first classified into when it entered the QFabric system). Traffic that matches the conditions in Step 5 is mapped back into its original forwarding class when it enters the egress Node device, after the traffic crosses the Interconnect device.

7. In the same firewall filter, configure a fourth term to define match conditions for traffic flowing from the Interconnect device to the egress Node device in the target-to-server direction (this term identifies traffic to map back into the original forwarding class after it crosses the Interconnect device):

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name from flow-match-conditions
user@switch# set filter filter-name term term-name from-fabric
```

8. In the fourth term in the same firewall filter, configure the action to take on traffic when it enters the egress Node device from the Interconnect device in the target-to-server direction (the action is to map the traffic back into its original forwarding class on the egress Node device, after the traffic crosses the Interconnect device):

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name then forwarding-class original-forwarding-class-name loss-priority priority-value
user@switch# set filter filter-name term term-name then counter counter-name
```

The first eight steps remap traffic in both directions of flow on the Interconnect device, and ensure that the traffic is mapped to its original forwarding class on the Node devices and as the traffic exits the QFabric system.

9. In the same firewall filter, add a final fifth term to define the default handling (action) of traffic that does not match the filter conditions. Firewall filters have an implicit default discard action, but in most cases, the intention is not to drop traffic is that is not remapped to a different forwarding class, so the action should be to accept the rest of the traffic:

```
[edit firewall family ethernet-switching]
user@switch# set filter filter-name term term-name then accept
```

If you wish, you can also configure a final counter in this term to count the total number of packets affected by the filter.

**NOTE:** If you configure a new forwarding class for the remapped traffic on a Node device, you must also configure scheduling for the new forwarding class on the Node device. On the Interconnect device, you must map the forwarding class to a fabric forwarding class set (fabric fc-set; see ["Understanding CoS Fabric Forwarding Class Sets" on page 396](#) for more information), and if the fabric fc-set is not one of the default fabric fc-sets, you must configure scheduling for the fabric fc-set (see ["Example: Configuring CoS Scheduling Across the QFabric System" on page 353](#) for more information).

## RELATED DOCUMENTATION

*Configuring Firewall Filters*

*Firewall Filter Match Conditions and Actions (QFX5100, QFX5110, QFX5120, QFX5200, EX4600, EX4650)*

[Example: Configuring CoS Scheduling Across the QFabric System | 353](#)

[Understanding How to Mitigate Fate Sharing on a QFabric System Interconnect Device by Remapping Traffic Flows \(Forwarding Classes\) | 410](#)

# 4

PART

## Configuring Data Center Bridging (ETS, PFC, DCBX) and Flow Control

---

Using Data Center Bridging and Flow Control | 443

---

# Using Data Center Bridging and Flow Control

## IN THIS CHAPTER

- [Understanding DCB Features and Requirements | 444](#)
- [Understanding CoS Hierarchical Port Scheduling \(ETS\) | 447](#)
- [Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 454](#)
- [Disabling the ETS Recommendation TLV | 490](#)
- [Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 491](#)
- [Enabling and Disabling CoS Symmetric Ethernet PAUSE Flow Control | 504](#)
- [Configuring CoS Asymmetric Ethernet PAUSE Flow Control | 505](#)
- [Configuring CoS PFC \(Congestion Notification Profiles\) | 507](#)
- [Example: Configuring CoS PFC for FCoE Traffic | 510](#)
- [Troubleshooting Dropped FCoE Traffic | 524](#)
- [Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 528](#)
- [Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic \(FCoE Transit Switch\) | 549](#)
- [Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface | 562](#)
- [Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces | 575](#)
- [Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications \(FCoE and iSCSI\) | 594](#)
- [Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway | 620](#)
- [Example: Configuring IEEE 802.1p Priority Remapping on an FCoE-FC Gateway | 624](#)
- [Understanding DCBX | 638](#)
- [Configuring the DCBX Mode | 648](#)
- [Configuring DCBX Autonegotiation | 649](#)
- [Understanding DCBX Application Protocol TLV Exchange | 652](#)
- [Defining an Application for DCBX Application Protocol TLV Exchange | 657](#)
- [Configuring an Application Map for DCBX Application Protocol TLV Exchange | 658](#)
- [Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange | 660](#)

- [Example: Configuring DCBX Application Protocol TLV Exchange | 661](#)

## Understanding DCB Features and Requirements

### IN THIS SECTION

- [Lossless Transport | 445](#)
- [ETS | 446](#)
- [DCBX | 447](#)

Data center bridging (DCB) is a set of enhancements to the IEEE 802.1 bridge specifications. DCB modifies and extends Ethernet behavior to support I/O convergence in the data center. I/O convergence includes but is not limited to the transport of Ethernet LAN traffic and Fibre Channel (FC) storage area network (SAN) traffic on the same physical Ethernet network infrastructure.



**Video:** [What is Data Center Bridging?](#)

A converged architecture saves cost by reducing the number of networks and switches required to support both types of traffic, reducing the number of interfaces required, reducing cable complexity, and reducing administration activities.

The Juniper Networks QFX Series and EX4600 switches support the DCB features required to transport converged Ethernet and FC traffic while providing the class-of-service (CoS) and other characteristics FC requires for transmitting storage traffic. To accommodate FC traffic, DCB specifications provide:

- A flow control mechanism called *priority-based flow control* (PFC, described in IEEE 802.1Qbb) to help provide lossless transport.
- A discovery and exchange protocol for conveying configuration and capabilities among neighbors to ensure consistent configuration across the network, called Data Center Bridging Capability Exchange protocol (DCBX), which is an extension of Link Layer Data Protocol (LLDP, described in IEEE 802.1AB).
- A bandwidth management mechanism called enhanced transmission selection (ETS, described in IEEE 802.1Qaz).

- A congestion management mechanism called quantized congestion notification (QCN, described in IEEE 802.1Qau).

The switch supports the PFC, DCBX, and ETS standards but does not support QCN. The switch also provides the high-bandwidth interfaces (10-Gbps minimum) required to support DCB and converged traffic.

This topic describes the DCB standards and requirements the switch supports:

## Lossless Transport

FC traffic requires lossless transport (defined as no frames dropped because of congestion). Standard Ethernet does not support lossless transport, but the DCB extensions to Ethernet along with proper buffer management enable an Ethernet network to provide the level of *class of service* (CoS) necessary to transport FC frames encapsulated in Ethernet over an Ethernet network.

This section describes these factors in creating lossless transport over Ethernet:

## PFC

PFC is a link-level flow control mechanism similar to Ethernet PAUSE (described in IEEE 802.3x). Ethernet PAUSE stops all traffic on a link for a period of time. PFC enables you to divide traffic on a link into eight priorities and stop the traffic of a selected priority without stopping the traffic assigned to other priorities on the link.

Pausing the traffic of a selected priority enables you to provide lossless transport for traffic assigned that priority and at the same time use standard lossy Ethernet transport for the rest of the link traffic.

## Buffer Management

Buffer management is critical to the proper functioning of PFC, because if buffers are allowed to overflow, frames are dropped and transport is not lossless.

For each lossless flow priority, the switch requires sufficient buffer space to:

- Store frames sent during the time it takes to send the PFC pause frame across the cable between devices.
- Store the frames that are already on the wire when the sender receives the PFC pause frame.

The propagation delay due to cable length and speed, as well as processing speed, determines the amount of buffer space needed to prevent frame loss due to congestion.

The switch automatically sets the threshold for sending PFC pause frames to accommodate delay from cables as long as 150 meters (492 feet) and to accommodate large frames that might be on the wire



when the switch sends the pause frame. This ensures that the switch sends pause frames early enough to allow the sender to stop transmitting before the receive buffers on the switch overflow.

## Physical Interfaces

QFX Series switches support 10-Gbps or faster, full-duplex interfaces. The switch enables DCB capability only on 10-Gbps or faster Ethernet interfaces.

## ETS

PFC divides traffic into up to eight separate streams (priorities, configured on the switch as forwarding classes) on a physical link. ETS enables you to manage the link bandwidth by:

- Grouping the priorities into priority groups (configured on the switch as forwarding class sets).
- Specifying the bandwidth available to each of the priority groups as a percentage of the total available link bandwidth.
- Allocating the bandwidth to the individual priorities in the priority group.

The available link bandwidth is the bandwidth remaining after servicing strict-high priority queues. On QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, we recommend that you always configure a shaping rate to limit the amount of bandwidth a strict-high priority queue can consume by including the [shaping-rate](#) statement in the [edit class-of-service schedulers] hierarchy on the strict-high priority scheduler. This prevents a strict-high priority queue from starving other queues on the port. (On QFX10000 switches, configure a transmit rate on strict-high priority queues to set a maximum amount of bandwidth for strict-high priority traffic.)

Managing link bandwidth with ETS provides several advantages:

- There is uniform management of all types of traffic on the link, both congestion-managed traffic and standard Ethernet traffic.
- When a priority group does not use all of its allocated bandwidth, other priority groups on the link can use that bandwidth as needed.

When a priority in a priority group does not use all of its allocated bandwidth, other priorities in the group can use that bandwidth.

The result is better bandwidth utilization, because priorities that consist of bursty traffic can share bandwidth during periods of low traffic transmission instead of consuming their entire bandwidth allocation when traffic loads are light.

- You can assign traffic types with different service needs to different priorities so that each traffic type receives appropriate treatment.

- Strict priority traffic retains its allocated bandwidth.

## DCBX

DCB devices use DCBX to exchange configuration information with directly connected peers (switches and endpoints such as servers). DCBX is an extension of LLDP. If you disable LLDP on an interface, that interface cannot run DCBX. If you attempt to enable DCBX on an interface on which LLDP is disabled, the configuration commit fails.

DCBX can:

- Discover the DCB capabilities of peers.
- Detect DCB feature misconfiguration or mismatches between peers.
- Configure DCB features on peers.

You can configure DCBX operation for PFC, ETS, and for Layer 2 and Layer 4 applications such as FCoE and iSCSI. DCBX is enabled or disabled on a per-interface basis.

## RELATED DOCUMENTATION

*Understanding FCoE*

[Understanding CoS Hierarchical Port Scheduling \(ETS\)](#)

*Understanding CoS Flow Control (Ethernet PAUSE and PFC)*

*Understanding DCBX*

*Example: Configuring CoS PFC for FCoE Traffic*

## Understanding CoS Hierarchical Port Scheduling (ETS)

### IN THIS SECTION

- [Hierarchical Scheduling Tiers | 448](#)
- [Hierarchical Scheduling and ETS | 449](#)
- [ETS Advertisement in DCBX | 451](#)
- [Hierarchical Scheduling Process | 451](#)
- [Strict-High Priority Queues and Hierarchical Scheduling | 453](#)

Scheduling defines the class-of-service (CoS) properties of output queues. Output queues are mapped to forwarding classes. CoS scheduler properties include the amount of interface bandwidth assigned to the queue, the queue priority, and the drop profiles associated with the queue.

Hierarchical port scheduling is a two-tier process that provides better port bandwidth utilization and greater flexibility to allocate resources to queues (forwarding classes) and to groups of queues (forwarding class sets). Hierarchical scheduling includes the Junos OS implementation of enhanced transmission selection (ETS), as described in IEEE 802.1Qaz.



Video: [What is Enhanced Transmission Selection?](#)

This topic describes:

**Hierarchical Scheduling Tiers**

The two tiers used in hierarchical scheduling are priorities and priority groups, as shown in [Table 84 on page 448](#).

**Table 84: Hierarchical Scheduling Tiers**

| Junos OS Configuration Construct | Equivalent ETS Construct | Description                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                        |
|----------------------------------|--------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Forwarding class                 | Priority                 | <p>Think about priorities (forwarding classes) as output queues. You map forwarding classes to queues, so each forwarding class represents an output queue.</p> <p>When you use a classifier to map a forwarding class to an IEEE 802.1p code point, the code point identifies that traffic's priority for priority-based flow control (PFC). Thus the forwarding class, the queue mapped to the forwarding class, and the priority (code point) mapped to the forwarding class all identify the same traffic.</p> |

**Table 84: Hierarchical Scheduling Tiers** *(Continued)*

| Junos OS Configuration Construct | Equivalent ETS Construct | Description                                                                                                                                                                                                                                                                                                   |
|----------------------------------|--------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Forwarding class set             | Priority group           | <p>Priority groups (forwarding class sets) are groups of priorities (forwarding classes). Forwarding class membership in a forwarding class set defines the priority group to which each priority belongs.</p> <p>You can configure up to three unicast priority groups and one multicast priority group.</p> |

You apply scheduling properties to each hierarchical scheduling tier as described in the next section.

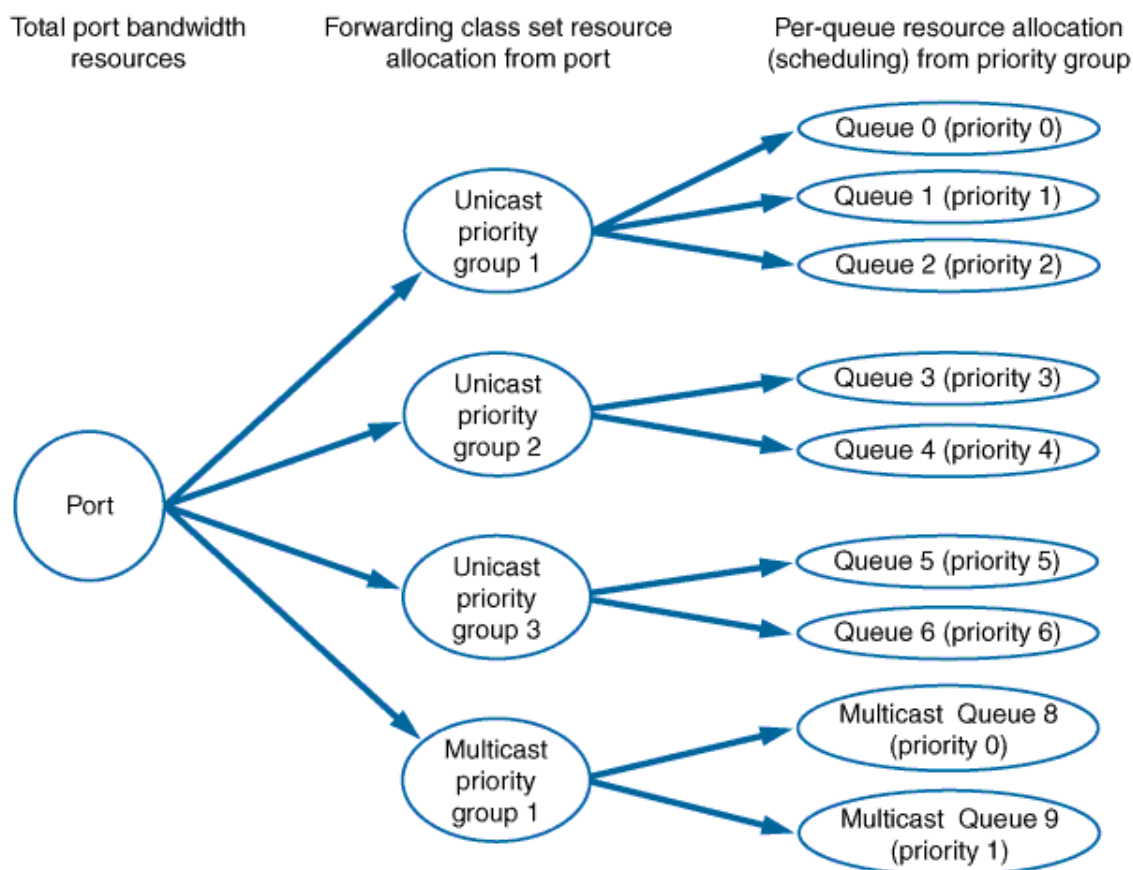
**NOTE:** If you explicitly configure one or more priority groups on an interface, any priority (forwarding class) that is not assigned to a priority group (forwarding class set) on that interface is assigned to an automatically generated default priority group and receives *no bandwidth*. This means that if you configure hierarchical scheduling on an interface, every forwarding class that you want to forward traffic on that interface must belong to a forwarding class set.

## Hierarchical Scheduling and ETS

Two-tier hierarchical scheduling manages bandwidth efficiently by enabling you to define the CoS properties for each priority group and for each priority. The first tier of the hierarchical scheduler allocates port bandwidth to a priority group. The second tier of the hierarchical scheduler determines the portion of the priority group bandwidth that a priority (queue) can use.

The CoS properties of a priority group define the amount of port bandwidth resources available to the queues in that priority group. The CoS properties you configure for each queue specify the amount of the bandwidth available to the queue from the bandwidth allocated to the priority group. [Figure 18 on page 450](#) shows the relationship of port resource allocation to priority groups, and priority group resource allocation to queues (priorities).

Figure 18: Hierarchical Scheduling Tiers



g040722

If a queue (priority) does not use its allocated bandwidth, ETS shares the unused bandwidth among the other queues in the priority group in proportion to the minimum guaranteed rate (transmit rate) scheduled for each queue. If a priority group does not use its allocated bandwidth, ETS shares the unused bandwidth among the priority groups on the port in proportion to the minimum guaranteed rate (guaranteed rate) scheduled for each priority group.

In this way, ETS improves link bandwidth utilization, and it provides each queue and each priority group with the maximum available bandwidth. For example, priorities that consist of bursty traffic can share bandwidth during periods of low traffic transmission, instead of reserving their entire bandwidth allocation when traffic loads are light.

**NOTE:** The available link bandwidth is the bandwidth remaining after servicing strict-high priority flows. Strict-high priority takes precedence over all other traffic. We recommend that you configure a *shaping-rate* (*transmit-rate* on QFX10000 switches) to limit the maximum amount of bandwidth that a strict-high priority forwarding class can use to prevent starving other queues.

## ETS Advertisement in DCBX

When you configure hierarchical scheduling on a port, Data Center Bridging Capability Exchange protocol (DCBX) advertises:

- Each priority group
- The priorities in each priority group
- The bandwidth properties of each priority group and priority

When you configure hierarchical scheduling on a port, any priority that is not part of an explicitly configured priority group is assigned to the automatically generated default priority group and receives no bandwidth. The default priority group is transparent. It does not appear in the configuration.

## Hierarchical Scheduling Process

Hierarchical scheduling consists of multiple configuration steps that create the priorities and the priority groups, schedule their resources, and assign them to interfaces. The steps below correspond to the six blocks in the packet flow diagram shown in [Figure 19 on page 452](#):

### 1. Packet classification:

- Configure classification of incoming traffic into forwarding classes (priorities). This consists of either using the default classifiers or configuring classifiers to map code points and loss priorities to the forwarding classes.
- Apply the classifiers to ingress interfaces or use the default classifiers. Applying a classifier to an interface groups incoming traffic on the interface into forwarding classes and loss priorities, by applying the classifier code point mapping to the incoming traffic.

### 2. Configure the output queues for the forwarding classes (priorities). This consists of either using the default forwarding classes and forwarding-class-to-queue mapping, or creating your own forwarding classes and mapping them to output queues.

### 3. Allocate resources to the forwarding classes:

- Define resources for the priorities. This consists of configuring schedulers to set minimum guaranteed bandwidth, maximum bandwidth, drop profiles for Weighted Random Early Detection (WRED), and bandwidth priority to apply to a forwarding class. Extra bandwidth is shared among queues in proportion to the minimum guaranteed bandwidth (transmit rate) of each queue.
- Map resources to priorities. This consists of mapping forwarding classes to schedulers, using a scheduler map.

### 4. Configure priority groups. This consists of mapping forwarding classes (priorities) to forwarding class sets (priority groups) to define the priorities that belong to each priority group.

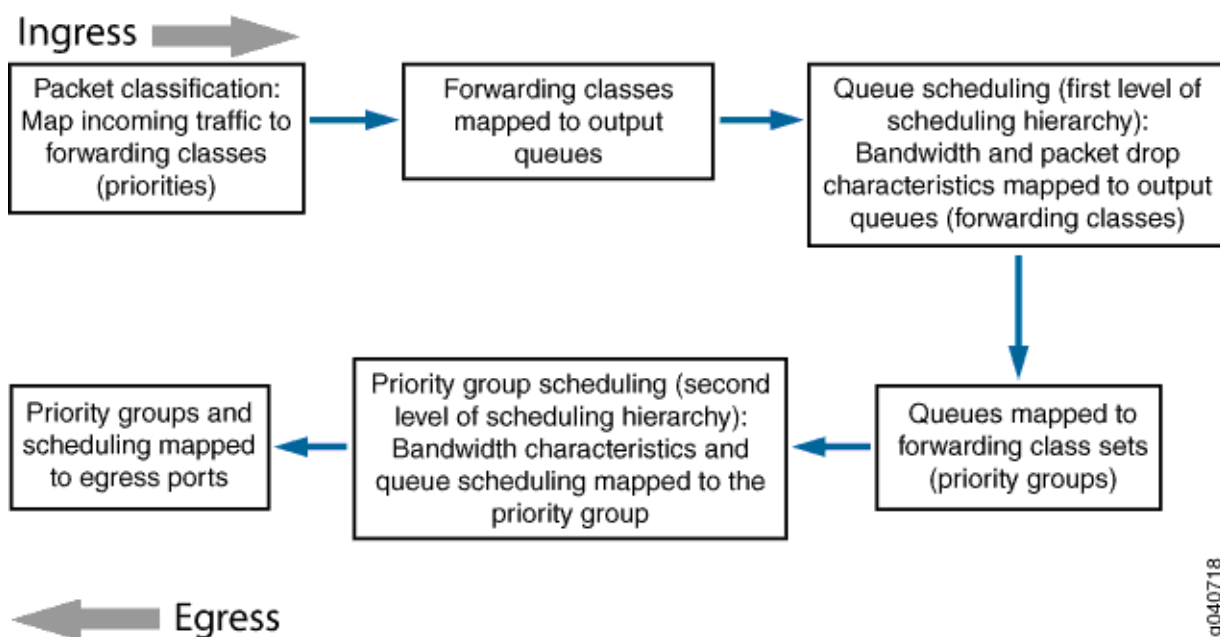
5. Define resources for the priority groups. This consists of configuring traffic control profiles to set minimum guaranteed bandwidth (*guaranteed-rate*) and maximum bandwidth (*shaping-rate* on switches other than QFX10000 switches, *transmit-rate* on QFX10000 switches) for a priority group. Traffic control profiles also specify a scheduler map, which defines the resources (schedulers) mapped to the priorities in the priority group. Extra port bandwidth is shared among priority groups in proportion to the minimum guaranteed bandwidth of each priority group.

The traffic control profile bandwidth settings determine the port resources available to the priority group. The schedulers specified in the scheduler map determine the amount of priority group resources that each priority receives.

**NOTE:** QFX10000 switches do not support defining a shaping rate for priority groups. Instead, set the maximum bandwidth for a priority group by defining a transmit rate. See *transmit-rate*.

6. Apply hierarchical scheduling to a port. This consists of attaching one or more priority groups (forwarding class sets) to an interface. For each priority group, you also attach a traffic control profile, which contains the scheduling properties of the priority group and the priorities in the priority group. Different priority groups on the same port can use different traffic control profiles, which provides fine tuned control of scheduling for each queue on each interface.

Figure 19: Hierarchical Scheduling Packet Flow



## Strict-High Priority Queues and Hierarchical Scheduling

If you configure a strict-high priority queue, you must observe the following rules:

- You must create a separate forwarding class set (priority group) for the strict-high priority queue.
- Only one forwarding class set can contain strict-high priority queues.
- Strict-high priority queues cannot belong to the same forwarding class set as queues that are not strict-high priority.
- A strict-high priority queue cannot belong to a multidestination forwarding class set.
- We recommend that you always apply a *shaping-rate* (*transmit-rate* on QFX10000 switches) to strict-high priority queues to limit the amount of bandwidth a strict-high priority queue can use. If you do not limit the amount of bandwidth a strict-high priority queue can use, then the strict-high priority queue can use all of the available port bandwidth and starve other queues on the port.

**NOTE:** On a QFabric system, if a fabric (fte) interface handles strict-high priority traffic, you must define a separate forwarding class set (priority group) for strict-high priority traffic. Strict-high priority traffic cannot be mixed with traffic of other priorities in a forwarding class set. For example, you might choose to create different forwarding class sets for best effort, lossless, strict-high priority, and multidestination traffic.

## Default Hierarchical Scheduling

**NOTE:** There is no default hierarchical scheduling on QFX10000 switches. QFX10000 switches use port scheduling by default, and you must explicitly configure hierarchical scheduling to enable ETS. Also on QFX10000 switches, changing from port scheduler to ETS or from ETS to port scheduler requires a reboot.

If you do not explicitly configure hierarchical scheduling, the switch uses the default settings:

- The switch automatically creates a default forwarding class set that contains all of the forwarding classes on the switch. The switch assigns 100 percent of the port output bandwidth to the default forwarding class set. The default forwarding class set is transparent. It does not appear in the configuration and is used for Data Center Bridging Capability Exchange protocol (DCBX) advertisement.
- Ingress traffic is classified based on the default classifier settings.



- The forwarding classes (queues) in the default forwarding class set receive bandwidth based on the default scheduler settings.

## RELATED DOCUMENTATION

*Understanding CoS Packet Flow*

*Understanding CoS Output Queue Schedulers*

*Understanding CoS Priority Group Scheduling*

[Benefits of Configuring CoS Hierarchical Port Scheduling](#)

*Understanding CoS Flow Control (Ethernet PAUSE and PFC)*

*Understanding CoS Classifiers*

*Understanding Default CoS Scheduling and Classification*

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

*Example: Configuring CoS Hierarchical Port Scheduling (ETS)*

*Example: Configuring Queue Schedulers*

*Example: Configuring Traffic Control Profiles (Priority Group Scheduling)*

*Example: Configuring Minimum Guaranteed Output Bandwidth*

*Example: Configuring Maximum Output Bandwidth*

## Example: Configuring CoS Hierarchical Port Scheduling (ETS)

### IN THIS SECTION

- [Requirements | 455](#)
- [Overview | 456](#)
- [Configuration | 462](#)
- [Verification | 476](#)

Hierarchical port scheduling defines the class-of-service (CoS) properties of output queues, which are mapped to forwarding classes. Traffic is classified into forwarding classes based on code point (priority), so mapping queues to forwarding classes also maps queues to priorities). Hierarchical port scheduling enables you to group priorities that require similar CoS treatment into priority groups. You define the port bandwidth resources for a priority group, and you define the amount of the priority group's resources that each priority in the group can use.

Hierarchical port scheduling is the Junos OS implementation of enhanced transmission selection (ETS), as described in IEEE 802.1Qaz. One major benefit of hierarchical port scheduling is greater port bandwidth utilization. If a priority group on a port does not use all of its allocated bandwidth, other priority groups on that port can use that bandwidth. Also, if a priority within a priority group does not use its allocated bandwidth, other priorities within that priority group can use that bandwidth.

Configuring hierarchical scheduling is a multistep procedure that includes:

- Mapping forwarding classes to queues
- Defining forwarding class sets (priority groups)
- Defining behavior aggregate classifiers
- Configuring priority-based flow control (PFC) for lossless priorities (queues)
- Applying classifiers and PFC configuration to ingress interfaces
- Defining drop profiles
- Defining schedulers
- Mapping forwarding classes to schedulers
- Defining traffic control profiles
- Assigning priority groups and traffic control profiles to egress ports

**NOTE:** OCX Series switches do not support lossless transport and do not support PFC. Although this example includes configuring lossless transport with PFC, the portions of the example that do not pertain to lossless transport still apply to OCX Series switches. (You can configure hierarchical scheduling on OCX Series switches, but you cannot configure lossless transport or lossless forwarding classes.)

This example describes how to configure hierarchical scheduling:

## Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 11.1 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

## Overview

### IN THIS SECTION

- [Topology | 457](#)

Keep the following considerations in mind when you plan the port bandwidth allocation for priority groups and for individual priorities:

- How much traffic and what types of traffic you expect to traverse the system.
- How you want to divide different types of traffic into priorities (forwarding classes) to apply different CoS treatment to different types of traffic. Dividing traffic into priorities includes:
  - Mapping the code points of ingress traffic to forwarding classes using behavior aggregate (BA) classifiers. This classifies incoming traffic into the appropriate forwarding class based on code point.
  - Mapping forwarding classes to output queues. This defines the output queue for each type of traffic.
  - Attaching the BA classifier to the desired ingress interfaces so that incoming traffic maps to the desired forwarding classes and queues.
- How you want to organize priorities into priority groups (forwarding class sets).

Traffic that requires similar treatment usually belongs in the same priority group. To do this, place forwarding classes that require similar bandwidth, loss, and other characteristics in the same forwarding class set. For example, you can map all types of best-effort traffic forwarding classes into one forwarding class set.

- How much of the port bandwidth you want to allocate to each priority group and to each of the priorities in each priority group. The following considerations apply to bandwidth allocation:
  - Estimate how much traffic you expect in each forwarding class, and how much traffic you expect in each forwarding class set (the amount of traffic you expect in a forwarding class set is the aggregate amount of traffic in the forwarding classes that belong to the forwarding class set).

- The combined minimum guaranteed bandwidth of the priorities (forwarding classes) in a priority group should not exceed the minimum guaranteed bandwidth of the priority group (forwarding class set). The transmit rate scheduler parameter defines the minimum guaranteed bandwidth for forwarding classes. Scheduler maps associate schedulers with forwarding classes.
- The combined minimum guaranteed bandwidth of the priority groups (forwarding class sets) on a port should not exceed the port's total bandwidth. The guaranteed rate parameter in the traffic control profile defines the minimum bandwidth for a forwarding class set. Associating a scheduler map with a traffic control profile sets the scheduling for the individual forwarding classes in the forwarding class set.

This example creates hierarchical port scheduling by defining priority groups for best effort, guaranteed delivery, and high-performance computing (HPC) traffic. Each priority group includes priorities that need to receive similar CoS treatment. Each priority group and each priority within each priority group receive the CoS resources needed to service their flows. Lossless priorities use PFC to prevent packet loss when the network experiences congestion.

### Topology

Table 62 on page 233 shows the configuration components for this example.

**NOTE:** OCX Series switches do not support lossless transport and do not support PFC. If you eliminate the configuration elements for the default lossless `fcoe` and `no-loss` forwarding classes (including classifier, forwarding class set, scheduler, and traffic control profile configuration for those forwarding classes) and for PFC, this example works for OCX Series switches. However, because the default `fcoe` and `no-loss` forwarding classes do not carry traffic on OCX Series switches, you can apply the bandwidth allocated to those forwarding classes to other forwarding classes. By default, the active forwarding classes (`best-effort`, `network-control`, and `mcast`) share the unused bandwidth assigned to the `fcoe` and `no-loss` forwarding classes.

**Table 85: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology**

| Property | Settings       |
|----------|----------------|
| Hardware | QFX3500 switch |

Table 85: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology (*Continued*)

| Property                                                                                                                 | Settings                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                  |
|--------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Mapping of forwarding classes (priorities) to queues                                                                     | <p>best-effort to queue 0</p> <p>be2 to queue 1</p> <p>fcoe (Fibre Channel over Ethernet) to queue 3</p> <p>no-loss to queue 4</p> <p>hpc (high-performance computing) to queue 5</p> <p>network-control to queue 7</p> <p><b>NOTE:</b> On switches that do not support the ELS CLI, if you are using Junos OS Release 12.2 or later, use the default forwarding-class-to-queue mapping for the lossless fcoe and no-loss forwarding classes. If you explicitly configure the default lossless forwarding classes, the traffic mapped to those forwarding classes is treated as lossy (best-effort) traffic and does <i>not</i> receive lossless treatment.</p> <p>On switches that do not support the ELS CLI, in Junos OS Release 12.3 and later, you can include the <i>no-loss</i> packet drop attribute in the explicit forwarding class configuration to configure a lossless forwarding class.</p> |
| Forwarding class sets (priority groups)                                                                                  | <p>best-effort-pg: contains forwarding classes best-effort, be2, and network control</p> <p>guar-delivery-pg: contains forwarding classes fcoe and no-loss</p> <p>hpc-pg: contains forwarding class hpc</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Behavior aggregate classifier (maps forwarding classes and loss priorities to incoming packets by IEEE 802.1 code point) | <p>Name—hsclassifier1</p> <p>Code point mapping:</p> <ul style="list-style-type: none"> <li>• 000 to forwarding class best-effort and loss priority low</li> <li>• 001 to forwarding class be2 and loss priority high</li> <li>• 011 to forwarding class fcoe and loss priority low</li> <li>• 100 to forwarding class no-loss and loss priority low</li> <li>• 101 to forwarding class hpc and loss priority low</li> <li>• 110 to forwarding class network-control and loss priority low</li> </ul>                                                                                                                                                                                                                                                                                                                                                                                                     |

Table 85: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology (*Continued*)

| Property                                                                                                                                       | Settings                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
|------------------------------------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| PFC                                                                                                                                            | <p>Congestion notification profile name—gd-cnp</p> <p>PFC enabled on code points: 011 (fcoe priority), 010 (no-loss priority)</p>                                                                                                                                                                                                                                                                                                                                                   |
| Drop profiles<br><br><b>NOTE:</b> The fcoe and no-loss priorities (queues) do not use drop profiles because they are lossless traffic classes. | <p>dp-be-low: drop start point 25, drop end point 50, maximum drop rate 80</p> <p>dp-be-high: drop start point 10, drop end point 40, maximum drop rate 100</p> <p>dp-hpc: drop start point 75, drop end point 90, maximum drop rate 75</p> <p>dp-nc: drop start point 80, drop end point 100, maximum drop rate 100</p>                                                                                                                                                            |
| Queue schedulers                                                                                                                               | <p>be-sched: minimum bandwidth 3g, maximum bandwidth 100%, priority low, drop profiles dp-be-low and dp-be-high</p> <p>fcoe-sched: minimum bandwidth 2.5g, maximum bandwidth 100%, priority low</p> <p>hpc-sched: minimum bandwidth 2g, maximum bandwidth 100%, priority low, drop profile dp-hpc</p> <p>nc-sched: minimum bandwidth 500m, maximum bandwidth 100%, priority low, drop profile dp-nc</p> <p>nl-sched: minimum bandwidth 2g, maximum bandwidth 100%, priority low</p> |
| Forwarding class-to-scheduler mapping                                                                                                          | <p>Scheduler map be-map:<br/>           Forwarding class best-effort, scheduler be-sched<br/>           Forwarding class be2, scheduler be-sched<br/>           Forwarding class network-control, scheduler nc-sched</p> <p>Scheduler map gd-map:<br/>           Forwarding class fcoe, scheduler fcoe-sched<br/>           Forwarding class no-loss, scheduler nl-sched</p> <p>Scheduler map hpc-map:<br/>           Forwarding class hpc, scheduler hpc-sched</p>                 |

**Table 85: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology (Continued)**

| Property                 | Settings                                                                                                                                                                                                                                                                                                                                                                                                                        |
|--------------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Traffic control profiles | <p>be-tcp: scheduler map be-map, minimum bandwidth 3.5g, maximum bandwidth 100%</p> <p>gd-tcp: scheduler map gd-map, minimum bandwidth 4.5g, maximum bandwidth 100%</p> <p>hpc-tcp: scheduler map hpc-map, minimum bandwidth 2g, maximum bandwidth 100%</p>                                                                                                                                                                     |
| Interfaces               | <p>This example configures hierarchical port scheduling on interfaces xe-0/0/20 and xe-0/0/21. Because traffic is bidirectional, you apply the ingress and egress configuration components to both interfaces:</p> <ul style="list-style-type: none"> <li>• Classifier Name—hsclassifier1</li> <li>• Forwarding class sets—best-effort-pg, guar-deliver-pg, hpc-pg</li> <li>• Congestion notification profile—gd-cnp</li> </ul> |

[Figure 7 on page 237](#) shows a block diagram of the configuration components and the configuration flow of the CLI statements used in the example. You can perform the configuration steps in a different sequence if you want.

Figure 20: Hierarchical Port Scheduling Components Block Diagram

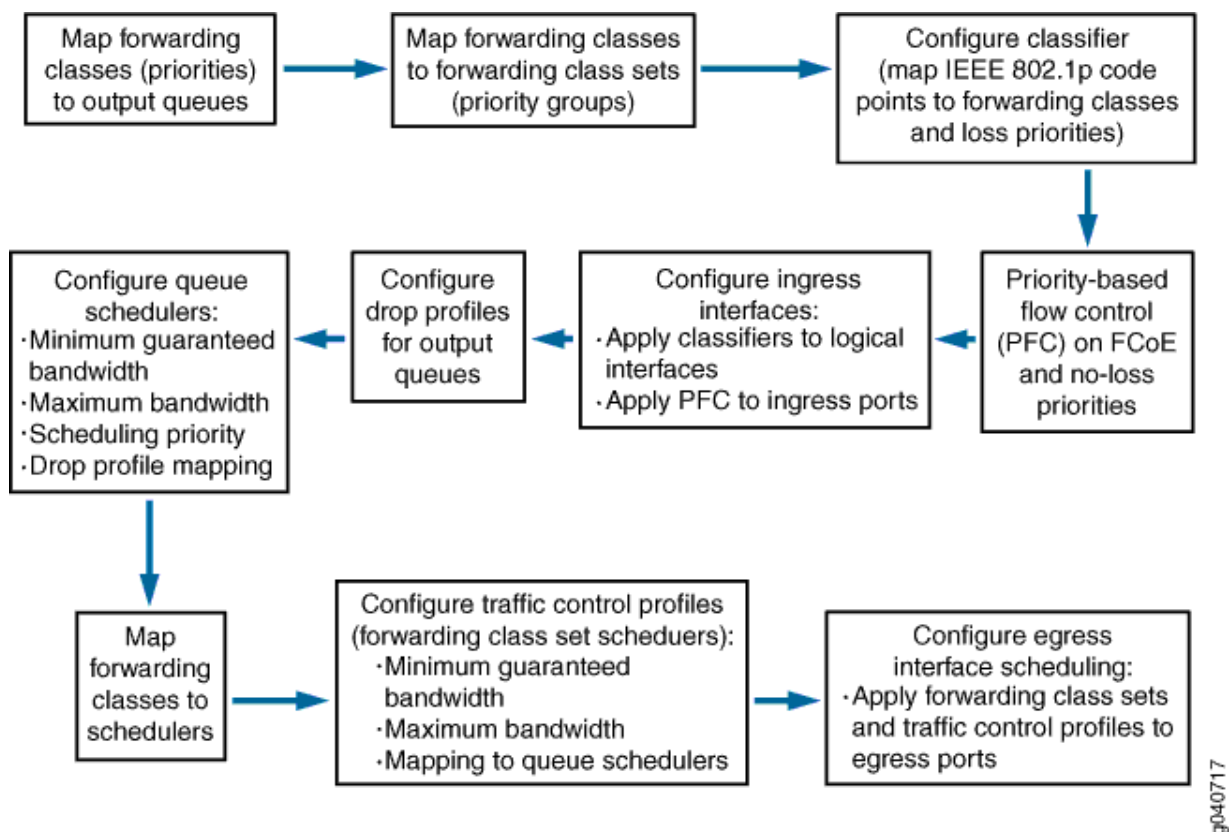
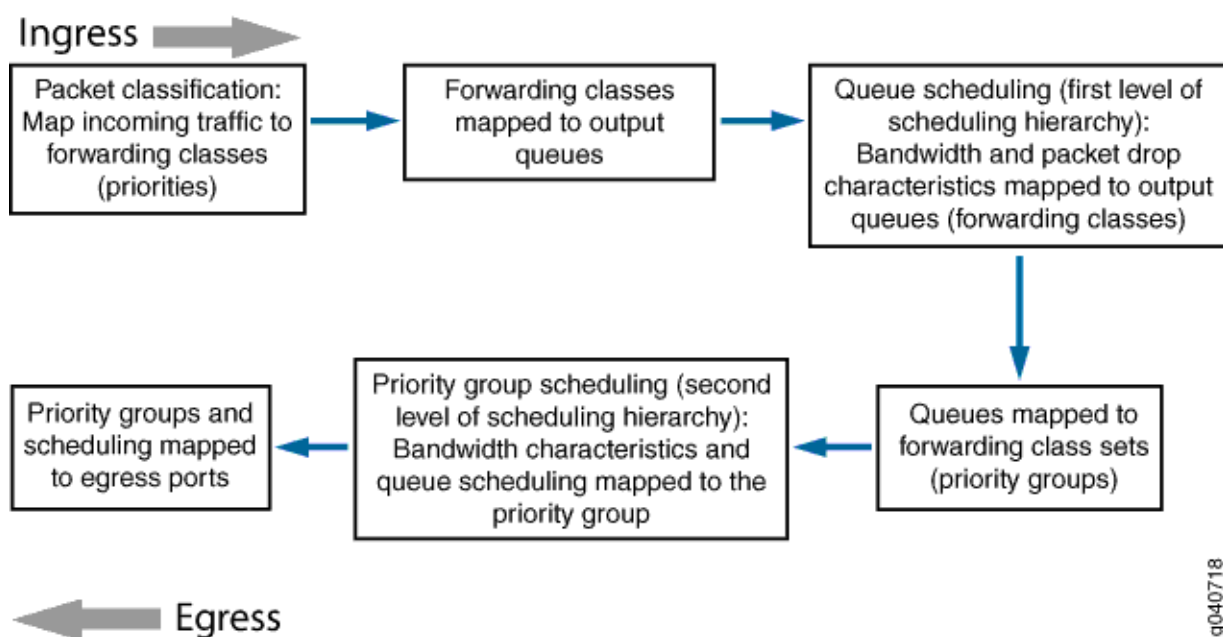


Figure 8 on page 238 shows a block diagram of the hierarchical scheduling packet flow from ingress to egress.



Figure 21: Hierarchical Port Scheduling Packet Flow Block Diagram



## Configuration

### IN THIS SECTION

- [CLI Quick Configuration | 462](#)
- [Procedure | 466](#)
- [Results | 472](#)

### CLI Quick Configuration

To quickly configure hierarchical port scheduling on systems that support lossless transport, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit class-of-service] hierarchy level:

```
[edit class-of-service]
set forwarding-classes class best-effort queue-num 0
set forwarding-classes class be2 queue-num 1
set forwarding-classes class hpc queue-num 5
```

```

set forwarding-classes class network-control queue-num 7
set forwarding-class-sets best-effort-pg class best-effort
set forwarding-class-sets best-effort-pg class be2
set forwarding-class-sets best-effort-pg class network-control
set forwarding-class-sets guar-delivery-pg class fcoe
set forwarding-class-sets guar-delivery-pg class no-loss
set forwarding-class-sets hpc-pg class hpc
set classifiers ieee-802.1 hsclassifier1 forwarding-class best-effort loss-priority low code-
points 000
set classifiers ieee-802.1 hsclassifier1 forwarding-class be2 loss-priority high code-points 001
set classifiers ieee-802.1 hsclassifier1 forwarding-class fcoe loss-priority low code-points
011
set classifiers ieee-802.1 hsclassifier1 forwarding-class no-loss loss-priority low code-points
100
set classifiers ieee-802.1 hsclassifier1 forwarding-class hpc loss-priority low code-points 101
set classifiers ieee-802.1 hsclassifier1 forwarding-class network-control loss-priority low code-
points 110
set congestion-notification-profile gd-cnp input ieee-802.1 code-point 011 pfc
set congestion-notification-profile gd-cnp input ieee-802.1 code-point 100 pfc
set interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 hsclassifier1
set interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 hsclassifier1
set interfaces xe-0/0/20 congestion-notification-profile gd-cnp
set interfaces xe-0/0/21 congestion-notification-profile gd-cnp
set drop-profiles dp-be-low interpolate fill-level 25 fill-level 50 drop-probability 0 drop-
probability 80
set drop-profiles dp-be-high interpolate fill-level 10 fill-level 40 drop-probability 0 drop-
probability 100
set drop-profiles dp-nc interpolate fill-level 80 fill-level 100 drop-probability 0 drop-
probability 100
set drop-profiles dp-hpc interpolate fill-level 75 fill-level 90 drop-probability 0 drop-
probability 75
set schedulers be-sched priority low transmit-rate 3g
set schedulers be-sched shaping-rate percent 100
set schedulers be-sched drop-profile-map loss-priority low protocol any drop-profile dp-be-low
set schedulers be-sched drop-profile-map loss-priority high protocol any drop-profile dp-be-high
set schedulers fcoe-sched priority low transmit-rate 2500m
set schedulers fcoe-sched shaping-rate percent 100
set schedulers hpc-sched priority low transmit-rate 2g
set schedulers hpc-sched shaping-rate percent 100
set schedulers hpc-sched drop-profile-map loss-priority low protocol any drop-profile dp-hpc
set schedulers nc-sched priority low transmit-rate 500m
set schedulers nc-sched shaping-rate percent 100
set schedulers nc-sched drop-profile-map loss-priority low protocol any drop-profile dp-nc

```

```

set schedulers nl-sched priority low transmit-rate 2g
set schedulers nl-sched shaping-rate percent 100
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set scheduler-maps be-map forwarding-class be2 scheduler be-sched
set scheduler-maps be-map forwarding-class network-control scheduler nc-sched
set scheduler-maps gd-map forwarding-class fcoe scheduler fcoe-sched
set scheduler-maps gd-map forwarding-class no-loss scheduler nl-sched
set scheduler-maps hpc-map forwarding-class hpc scheduler hpc-sched
set traffic-control-profiles be-tcp scheduler-map be-map guaranteed-rate 3500m
set traffic-control-profiles be-tcp shaping-rate percent 100
set traffic-control-profiles gd-tcp scheduler-map gd-map guaranteed-rate 4500m
set traffic-control-profiles gd-tcp shaping-rate percent 100
set traffic-control-profiles hpc-tcp scheduler-map hpc-map guaranteed-rate 2g
set traffic-control-profiles hpc-tcp shaping-rate percent 100
set interfaces xe-0/0/20 forwarding-class-set best-effort-pg output-traffic-control-profile be-
tcp
set interfaces xe-0/0/20 forwarding-class-set guar-delivery-pg output-traffic-control-profile gd-
tcp
set interfaces xe-0/0/20 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp
set interfaces xe-0/0/21 forwarding-class-set best-effort-pg output-traffic-control-profile be-
tcp
set interfaces xe-0/0/21 forwarding-class-set guar-delivery-pg output-traffic-control-profile gd-
tcp
set interfaces xe-0/0/21 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp

```

## OCX Series Switches

Because OCX Series switches do not support lossless transport, the following subset of the configuration eliminates the lossless configuration elements and provides hierarchical port scheduling for the best-effort, be2, hpc, and network-control forwarding classes. In addition, on OCX Series switches, you would probably use DSCP classifiers and code points instead of IEEE classifiers and code points. To quickly configure hierarchical port scheduling on an OCX Series switch, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit class-of-service] hierarchy level:

```

[edit class-of-service]
set forwarding-classes class best-effort queue-num 0
set forwarding-classes class be2 queue-num 1
set forwarding-classes class hpc queue-num 5
set forwarding-classes class network-control queue-num 7
set forwarding-class-sets best-effort-pg class best-effort

```

```

set forwarding-class-sets best-effort-pg class be2
set forwarding-class-sets best-effort-pg class network-control

set forwarding-class-sets hpc-pg class hpc
set classifiers ieee-802.1 hsclassifier1 forwarding-class best-effort loss-priority low code-points 000
set classifiers ieee-802.1 hsclassifier1 forwarding-class be2 loss-priority high code-points 001

set classifiers ieee-802.1 hsclassifier1 forwarding-class hpc loss-priority low code-points 101
set classifiers ieee-802.1 hsclassifier1 forwarding-class network-control loss-priority low code-points 110

set interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 hsclassifier1
set interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 hsclassifier1
set drop-profiles dp-be-low interpolate fill-level 25 fill-level 50 drop-probability 0 drop-probability 80
set drop-profiles dp-be-high interpolate fill-level 10 fill-level 40 drop-probability 0 drop-probability 100
set drop-profiles dp-nc interpolate fill-level 80 fill-level 100 drop-probability 0 drop-probability 100
set drop-profiles dp-hpc interpolate fill-level 75 fill-level 90 drop-probability 0 drop-probability 75
set schedulers be-sched priority low transmit-rate 3g
set schedulers be-sched shaping-rate percent 100
set schedulers be-sched drop-profile-map loss-priority low protocol any drop-profile dp-be-low
set schedulers be-sched drop-profile-map loss-priority high protocol any drop-profile dp-be-high
set schedulers hpc-sched priority low transmit-rate 2g
set schedulers hpc-sched shaping-rate percent 100
set schedulers hpc-sched drop-profile-map loss-priority low protocol any drop-profile dp-hpc
set schedulers nc-sched priority low transmit-rate 500m
set schedulers nc-sched shaping-rate percent 100
set schedulers nc-sched drop-profile-map loss-priority low protocol any drop-profile dp-nc
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set scheduler-maps be-map forwarding-class be2 scheduler be-sched
set scheduler-maps be-map forwarding-class network-control scheduler nc-sched
set scheduler-maps hpc-map forwarding-class hpc scheduler hpc-sched
set traffic-control-profiles be-tcp scheduler-map be-map guaranteed-rate 3500m
set traffic-control-profiles be-tcp shaping-rate percent 100
set traffic-control-profiles hpc-tcp scheduler-map hpc-map guaranteed-rate 2g
set traffic-control-profiles hpc-tcp shaping-rate percent 100
set interfaces xe-0/0/20 forwarding-class-set best-effort-pg output-traffic-control-profile be-tcp
set interfaces xe-0/0/20 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp

```

```
set interfaces xe-0/0/21 forwarding-class-set best-effort-pg output-traffic-control-profile be-
tcp
set interfaces xe-0/0/21 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp
```

## Procedure

### Step-by-Step Procedure

To perform a step-by-step configuration of the forwarding classes (priorities), forwarding class sets (priority groups), classifiers, queue schedulers, PFC, traffic control profiles, and interfaces to set up hierarchical port scheduling (ETS):

1. Configure the forwarding classes (priorities) and map them to unicast output queues (do not explicitly map the `fcoe` and `no-loss` forwarding classes to output queues; use the default configuration):

```
[edit class-of-service]
user@switch# set forwarding-classes class best-effort queue-num 0
user@switch# set forwarding-classes class be2 queue-num 1
user@switch# set forwarding-classes class hpc queue-num 5
user@switch# set forwarding-classes class network-control queue-num 7
```

2. Configure forwarding class sets (priority groups) to group forwarding classes (priorities) that require similar CoS treatment:

```
[edit class-of-service]
user@switch# set forwarding-class-sets best-effort-pg class best-effort
user@switch# set forwarding-class-sets best-effort-pg class be2
user@switch# set forwarding-class-sets best-effort-pg class network-control
user@switch# set forwarding-class-sets guar-delivery-pg class fcoe
user@switch# set forwarding-class-sets guar-delivery-pg class no-loss
user@switch# set forwarding-class-sets hpc-pg class hpc
```

**NOTE:** On OCX Series switches, you would not configure the `guar-delivery-pg` forwarding class set for lossless traffic.

3. Configure a classifier to set the loss priority and IEEE 802.1 code points assigned to each forwarding class at the ingress:

```
[edit class-of-service]
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class best-effort loss-
priority low code-points 000
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class be2 loss-priority
high code-points 001
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class fcoe loss-priority
low code-points 011
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class no-loss loss-
priority low code-points 100
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class hpc loss-priority
low code-points 101
user@switch# set classifiers ieee-802.1 hsclassifier1 forwarding-class network-control loss-
priority low code-points 110
```

**NOTE:** On OCX Series switches, you would not configure the fcoe and no-loss portions of the classifier.

4. Configure a congestion notification profile to enable PFC on the FCoE and no-loss queue IEEE 802.1 code points:

```
[edit class-of-service]
user@switch# set congestion-notification-profile gd-cnp input ieee-802.1 code-point 011 pfc
user@switch# set congestion-notification-profile gd-cnp input ieee-802.1 code-point 100 pfc
```

**NOTE:** This step does not apply to OCX Series switches, which do not support PFC.

5. Assign the classifier to the interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 hsclassifier1
user@switch# set interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 hsclassifier1
```

6. Apply the PFC configuration to the interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 congestion-notification-profile gd-cnp
user@switch# set interfaces xe-0/0/21 congestion-notification-profile gd-cnp
```

**NOTE:** This step does not apply to OCX Series switches, which do not support PFC.

7. Configure the drop profile for the best-effort low loss-priority queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-be-low interpolate fill-level 25 fill-level 50 drop-
probability 0 drop-probability 80
```

8. Configure the drop profile for the best-effort high loss-priority queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-be-high interpolate fill-level 10 fill-level 40 drop-
probability 0 drop-probability 100
```

9. Configure the drop profile for the network-control queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-nc interpolate fill-level 80 fill-level 100 drop-
probability 0 drop-probability 100
```

10. Configure the drop profile for the high-performance computing queue:

```
[edit class-of-service]
user@switch# set drop-profiles dp-hpc interpolate fill-level 75 fill-level 90 drop-
probability 0 drop-probability 75
```

11. Define the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profiles for the best-effort queue:

```
[edit class-of-service]
user@switch# set schedulers be-sched priority low transmit-rate 3g
user@switch# set schedulers be-sched shaping-rate percent 100
user@switch# set schedulers be-sched drop-profile-map loss-priority low protocol any drop-
profile dp-be-low
user@switch# set schedulers be-sched drop-profile-map loss-priority high protocol any drop-
profile dp-be-high
```

12. Define the minimum guaranteed bandwidth, priority, and maximum bandwidth for the FCoE queue:

```
[edit class-of-service]
user@switch# set schedulers fcoe-sched priority low transmit-rate 2500m
user@switch# set schedulers fcoe-sched shaping-rate percent 100
```

**NOTE:** This step does not apply to OCX Series switches, which do not support lossless transport.

13. Define the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profile for the high-performance computing queue:

```
[edit class-of-service]
user@switch# set schedulers hpc-sched priority low transmit-rate 2g
user@switch# set schedulers hpc-sched shaping-rate percent 100
user@switch# set schedulers hpc-sched drop-profile-map loss-priority low protocol any drop-
profile dp-hpc
```

14. Define the minimum guaranteed bandwidth, priority, maximum bandwidth, and drop profile for the network-control queue:

```
[edit class-of-service]
user@switch# set schedulers nc-sched priority low transmit-rate 500m
user@switch# set schedulers nc-sched shaping-rate percent 100
```



```
user@switch# set schedulers nc-sched drop-profile-map loss-priority low protocol any drop-profile dp-nc
```

15. Define the minimum guaranteed bandwidth, priority, and maximum bandwidth for the no-loss queue:

```
[edit class-of-service]
user@switch# set schedulers nl-sched priority low transmit-rate 2g
user@switch# set schedulers nl-sched shaping-rate percent 100
```

**NOTE:** This step does not apply to OCX Series switches, which do not support lossless transport.

16. Map the schedulers to the appropriate forwarding classes (queues):

```
[edit class-of-service]
user@switch# set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
user@switch# set scheduler-maps be-map forwarding-class be2 scheduler be-sched
user@switch# set scheduler-maps be-map forwarding-class network-control scheduler nc-sched
user@switch# set scheduler-maps gd-map forwarding-class fcoe scheduler fcoe-sched
user@switch# set scheduler-maps gd-map forwarding-class no-loss scheduler nl-sched
user@switch# set scheduler-maps hpc-map forwarding-class hpc scheduler hpc-sched
```

**NOTE:** On OCX Series switches, because lossless transport is not supported, you would not configure the `gd-map` scheduler map.

17. Define the traffic control profile for the best-effort priority group (queue scheduler to mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
user@switch# set traffic-control-profiles be-tcp scheduler-map be-map guaranteed-rate 3500m
user@switch# set traffic-control-profiles be-tcp shaping-rate percent 100
```

18. Define the traffic control profile for the guaranteed delivery priority group (queue to scheduler mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
user@switch# set traffic-control-profiles gd-tcp scheduler-map gd-map guaranteed-rate 4500m
user@switch# set traffic-control-profiles gd-tcp shaping-rate percent 100
```

**NOTE:** This step does not apply to OCX Series switches, which do not support lossless transport.

19. Define the traffic control profile for the high-performance computing priority group (queue to scheduler mapping, minimum guaranteed bandwidth, and maximum bandwidth):

```
[edit class-of-service]
user@switch# set traffic-control-profiles hpc-tcp scheduler-map hpc-map guaranteed-rate 2g
user@switch# set traffic-control-profiles hpc-tcp shaping-rate percent 100
```

20. Apply the three priority groups (forwarding class sets) and the appropriate traffic control profiles to the egress ports:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 forwarding-class-set best-effort-pg output-traffic-control-profile be-tcp
user@switch# set interfaces xe-0/0/20 forwarding-class-set guar-delivery-pg output-traffic-control-profile gd-tcp
user@switch# set interfaces xe-0/0/20 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp
user@switch# set interfaces xe-0/0/21 forwarding-class-set best-effort-pg output-traffic-control-profile be-tcp
user@switch# set interfaces xe-0/0/21 forwarding-class-set guar-delivery-pg output-traffic-control-profile gd-tcp
user@switch# set interfaces xe-0/0/21 forwarding-class-set hpc-pg output-traffic-control-profile hpc-tcp
```

**NOTE:** Because OCX Series switches do not support lossless transport, on OCX Series switches, you would not apply the `guar-deliver-pg` forwarding class set and the `gd-tcp` traffic control profile to interfaces.

## Results

Display the results of the configuration (the system shows only the explicitly configured parameters; it does not show default parameters such as the `fcoe` and `no-loss` lossless forwarding classes). On OCX Series switches, you would not see the lossless configuration components in the output:

```
user@switch> show configuration class-of-service
classifiers {
 ieee-802.1 hsclassifier1 {
 forwarding-class best-effort {
 loss-priority low code-points 000;
 }
 forwarding-class be2 {
 loss-priority high code-points 001;
 }
 forwarding-class fcoe {
 loss-priority low code-points 011;
 }
 forwarding-class no-loss {
 loss-priority low code-points 100;
 }
 forwarding-class hpc {
 loss-priority low code-points 101;
 }
 forwarding-class network-control {
 loss-priority low code-points 110;
 }
 }
}
drop-profiles {
 dp-be-low {
 interpolate {
 fill-level [25 50];
 drop-probability [0 80];
 }
 }
}
```

```

dp-be-high {
 interpolate {
 fill-level [10 40];
 drop-probability [0 100];
 }
}
dp-hpc {
 interpolate {
 fill-level [75 90];
 drop-probability [0 75];
 }
}
dp-nc {
 interpolate {
 fill-level [80 100];
 drop-probability [0 100];
 }
}
}
forwarding-classes {
 class best-effort queue-num 0;
 class be2 queue-num 1;
 class hpc queue-num 5;
 class network-control queue-num 7;
}
traffic-control-profiles {
 be-tcp {
 scheduler-map be-map;
 shaping-rate percent 100;
 guaranteed-rate 3500000000;
 }
 gd-tcp {
 scheduler-map gd-map;
 shaping-rate percent 100;
 guaranteed-rate 4500000000;
 }
 hpc-tcp {
 scheduler-map hpc-map;
 shaping-rate percent 100;
 guaranteed-rate 2g;
 }
}
forwarding-class-sets {

```

```

 guar-delivery-pg {
 class fcoe;
 class no-loss;
 }
 best-effort-pg {
 class best-effort;
 class be2;
 class network-control;
 }
 hpc-pg {
 class hpc;
 }
}
congestion-notification-profile {
 gd-cnp {
 input {
 ieee-802.1 {
 code-point 011 {
 pfc;
 }
 code-point 100 {
 pfc;
 }
 }
 }
 }
}
}
interfaces {
 xe-0/0/20 {
 forwarding-class-set {
 best-effort-pg {
 output-traffic-control-profile be-tcp;
 }
 guar-delivery-pg {
 output-traffic-control-profile gd-tcp;
 }
 hpc-pg {
 output-traffic-control-profile hpc-tcp;
 }
 }
 congestion-notification-profile gd-cnp;
 unit 0 {
 classifiers {

```

```

 ieee-802.1 hsclassifier1;
 }
}
xe-0/0/21 {
 forwarding-class-set {
 best-effort-pg {
 output-traffic-control-profile be-tcp;
 }
 guar-delivery-pg {
 output-traffic-control-profile gd-tcp;
 }
 hpc-pg {
 output-traffic-control-profile hpc-tcp;
 }
 }
 congestion-notification-profile gd-cnp;
 unit 0 {
 classifiers {
 ieee-802.1 hsclassifier1;
 }
 }
}
scheduler-maps {
 be-map {
 forwarding-class best-effort scheduler be-sched;
 forwarding-class network-control scheduler nc-sched;
 forwarding-class be2 scheduler be-sched;
 }
 gd-map {
 forwarding-class fcoe scheduler fcoe-sched;
 forwarding-class no-loss scheduler nl-sched;
 }
 hpc-map {
 forwarding-class hpc scheduler hpc-sched;
 }
}
schedulers {
 be-sched {
 transmit-rate 3g;
 shaping-rate percent 100;
 priority low;
 }
}

```

```

 drop-profile-map loss-priority low protocol any drop-profile dp-be-low;
 drop-profile-map loss-priority high protocol any drop-profile dp-be-high;
 }
 fcoe-sched {
 transmit-rate 2500000000;
 shaping-rate percent 100;
 priority low;
 }
 hpc-sched {
 transmit-rate 2g;
 shaping-rate percent 100;
 priority low;
 drop-profile-map loss-priority low protocol any drop-profile dp-hpc;
 }
 nc-sched {
 transmit-rate 500m;
 shaping-rate percent 100;
 priority low;
 drop-profile-map loss-priority low protocol any drop-profile dp-nc;
 }
 nl-sched {
 transmit-rate 2g;
 shaping-rate percent 100;
 priority low;
 }
}

```

**TIP:** To quickly configure the interfaces, issue the `load merge` terminal command, and then copy the hierarchy and paste it into the switch terminal window.

## Verification

### IN THIS SECTION

- [Verifying the Forwarding Classes \(Priorities\) | 477](#)
- [Verifying the Forwarding Class Sets \(Priority Groups\) | 478](#)
- [Verifying the Classifier | 479](#)
- [Verifying Priority-Based Flow Control | 480](#)

- [Verifying the Output Queue Schedulers | 481](#)
- [Verifying the Drop Profiles | 485](#)
- [Verifying the Priority Group Output Schedulers \(Traffic Control Profiles\) | 486](#)
- [Verifying the Interface Configuration | 487](#)

**NOTE:** The verification output is based on the full example configuration. On OCX Series switches, you do not see lossless configuration components in the output. Comments about lossless configuration components do not apply to OCX Series switches.

To verify that you created the hierarchical port scheduling components and they are operating properly, perform these tasks:

### Verifying the Forwarding Classes (Priorities)

#### Purpose

Verify that you created the forwarding classes and mapped them to the correct queues. (The system shows only the explicitly configured forwarding classes. It does not show default forwarding classes such as fcoe and no-loss.)

#### Action

List the forwarding classes using the operational mode command `show class-of-service forwarding-class`:

```
user@switch> show class-of-service forwarding-class
```

| Forwarding class | ID | Queue | Policing priority | No-Loss  |
|------------------|----|-------|-------------------|----------|
| best-effort      | 0  | 0     | normal            | Disabled |
| be2              | 1  | 3     | normal            | Disabled |
| hpc              | 2  | 4     | normal            | Disabled |
| network-control  | 3  | 7     | normal            | Disabled |
| mcast            | 8  | 8     | normal            | Disabled |



## Meaning

The `show class-of-service forwarding-class` command lists all of the configured forwarding classes, the internal identification number of each forwarding class, the queues that are mapped to the forwarding classes, the policing priority, and whether the forwarding class is lossless (no-loss packet drop attribute enabled) or lossy forwarding class (no-loss packet drop attribute disabled). The command output shows that:

- Forwarding class `best-effort` maps to queue 0 and is lossy
- Forwarding class `be2` maps to queue 1 and is lossy
- Forwarding class `hpc` maps to queue 5 and is lossy
- Forwarding class `network-control` maps to queue 7 and is lossy

In addition, the command lists the default multicast (multidestination) forwarding class and the default queue to which it is mapped.

## Verifying the Forwarding Class Sets (Priority Groups)

### Purpose

Verify that you created the priority groups and that the correct priorities (forwarding classes) belong to the appropriate priority group.

### Action

List the forwarding class sets using the operational mode command `show class-of-service forwarding-class-set`:

```
user@switch> show class-of-service forwarding-class-set
Forwarding class set: best-effort-pg, Type: normal-type, Forwarding class set index: 19907
 Forwarding class Index
 best-effort 0
 be2 1
 network-control 5

Forwarding class set: guar-delivery-pg, Type: normal-type, Forwarding class set index: 43700
 Forwarding class Index
 fcoe 2
 no-loss 3
```

```
Forwarding class set: hpc-pg, Type: normal-type, Forwarding class set index: 60758
```

| Forwarding class | Index |
|------------------|-------|
| hpc              | 4     |

## Meaning

The `show class-of-service forwarding-class-set` command lists all of the configured forwarding class sets (priority groups), the forwarding classes (priorities) that belong to each priority group, and the internal index number of each priority group. The command output shows that:

- The forwarding class set `best-effort-pg` includes the forwarding classes `best-effort`, `be2`, and `network-control`.
- The forwarding class set `guar-delivery-pg` includes the forwarding classes `fcoe` and `no-loss`.
- The forwarding class set `hpc-pg` includes the forwarding class `hpc`.

## Verifying the Classifier

### Purpose

Verify that the classifier maps forwarding classes to the correct IEEE 802.1p code points and packet loss priorities.

### Action

List the classifier configured for hierarchical port scheduling using the operational mode command `show class-of-service classifier name hsclassifier1`:

```
user@switch> show class-of-service classifier name hsclassifier1
Classifier: hsclassifier1, Code point type: ieee-802.1, Index: 43607
 Code point Forwarding class Loss priority
 000 best-effort low
 001 be2 high
 011 fcoe low
 100 no-loss low
 101 hpc low
 110 network-control low
```

### Meaning

The `show class-of-service classifier name hsclassifier1` command lists all of the IEEE 802.1p code points and the loss priorities mapped to all of the forwarding classes in the classifier. The command output shows that the forwarding classes `best-effort`, `be2`, `no-loss`, `fcoe`, `hpc`, and `network-control` have been created and mapped to IEEE 802.1p code points and loss priorities.

### Verifying Priority-Based Flow Control

#### Purpose

Verify that PFC is enabled on the correct priorities for lossless transport.

#### Action

List the congestion notification profiles using the operational mode command `show class-of-service congestion-notification`:

```

user@switch> show class-of-service congestion-notification
Type: Input, Name: gd-cnp, Index: 51687
Cable Length: 100 m
 Priority PFC MRU
 000 Disabled
 001 Disabled
 010 Disabled
 011 Enabled 2500
 100 Enabled 2500
 101 Disabled
 110 Disabled
 111 Disabled
Type: Output
 Priority Flow-Control-Queues
 000
 001 0
 010 1
 011 2
 011 3
 100

```

|     |   |
|-----|---|
|     | 4 |
| 101 |   |
|     | 5 |
| 110 |   |
|     | 6 |
| 111 |   |
|     | 7 |

## Meaning

The `show class-of-service congestion-notification` command lists all of the congestion notification profiles and the IEEE 802.1p code points with PFC enabled. The command output shows that PFC is enabled for code points 011 (fcoe priority and queue) and 100 (no-loss priority and queue) for the `gd-cnp` congestion notification profile.

The command also shows the default cable length (100 meters), the default maximum receive unit (2500 bytes), and the default mapping of priorities to output queues because this example does not include configuring these options.

## Verifying the Output Queue Schedulers

### Purpose

Verify that you created the output queue schedulers with the correct bandwidth parameters and priorities, mapped to the correct queues, and mapped to the correct drop profiles.

### Action

List the scheduler maps using the operational mode command `show class-of-service scheduler-map`:

```
user@switch> show class-of-service scheduler-map
Scheduler map: be-map, Index: 64023

Scheduler: be-sched, Forwarding class: best-effort, Index: 13005
 Transmit rate: 3000000000 bps, Rate Limit: none, Buffer size: remainder,
 Buffer Limit: none, Priority: low
 Excess Priority: unspecified
 Shaping rate: 100 percent,
 drop-profile-map-set-type: mark
 Drop profiles:
 Loss priority Protocol Index Name
```

|             |     |       |                        |
|-------------|-----|-------|------------------------|
| Low         | any | 55387 | dp-be-low              |
| Medium high | any | 1     | <default-drop-profile> |
| High        | any | 4369  | dp-be-high             |

Scheduler: be-sched, Forwarding class: be2, Index: 13005

Transmit rate: 3000000000 bps, Rate Limit: none, Buffer size: remainder,

Buffer Limit: none, Priority: low

Excess Priority: unspecified

Shaping rate: 100 percent,

drop-profile-map-set-type: mark

Drop profiles:

| Loss priority | Protocol | Index | Name                   |
|---------------|----------|-------|------------------------|
| Low           | any      | 55387 | dp-be-low              |
| Medium high   | any      | 1     | <default-drop-profile> |
| High          | any      | 4369  | dp-be-high             |

Scheduler: nc-sched, Forwarding class: network-control, Index: 45740

Transmit rate: 5000000000 bps, Rate Limit: none, Buffer size: remainder,

Buffer Limit: none, Priority: low

Excess Priority: unspecified

Shaping rate: 100 percent,

drop-profile-map-set-type: mark

Drop profiles:

| Loss priority | Protocol | Index | Name                   |
|---------------|----------|-------|------------------------|
| Low           | any      | 44207 | dp-nc                  |
| Medium high   | any      | 1     | <default-drop-profile> |
| High          | any      | 1     | <default-drop-profile> |

Scheduler map: gd-map, Index: 61447

Scheduler: fcoe-sched, Forwarding class: fcoe, Index: 37289

Transmit rate: 2500000000 bps, Rate Limit: none, Buffer size: remainder,

Buffer Limit: none, Priority: low

Excess Priority: unspecified

Shaping rate: 100 percent,

drop-profile-map-set-type: mark

Drop profiles:

| Loss priority | Protocol | Index | Name                   |
|---------------|----------|-------|------------------------|
| Low           | any      | 44207 | <default-drop-profile> |
| Medium high   | any      | 1     | <default-drop-profile> |
| High          | any      | 1     | <default-drop-profile> |

Scheduler: nl-sched, Forwarding class: no-loss, Index: 29359

```

Transmit rate: 2000000000 bps, Rate Limit: none, Buffer size: remainder,
Buffer Limit: none, Priority: low
Excess Priority: unspecified
Shaping rate: 100 percent,
drop-profile-map-set-type: mark
Drop profiles:
 Loss priority Protocol Index Name
 Low any 44207 <default-drop-profile>
 Medium high any 1 <default-drop-profile>
 High any 1 <default-drop-profile>

```

Scheduler map: hpc-map, Index: 56941

```

Scheduler: hpc-sched, Forwarding class: hpc, Index: 55900
Transmit rate: 2000000000 bps, Rate Limit: none, Buffer size: remainder,
Buffer Limit: none, Priority: low
Excess Priority: unspecified
Shaping rate: 100 percent,
drop-profile-map-set-type: mark
Drop profiles:
 Loss priority Protocol Index Name
 Low any 57716 dp-hpc
 Medium high any 1 <default-drop-profile>
 High any 1 <default-drop-profile>

```

## Meaning

The `show class-of-service scheduler-map` command lists all of the configured scheduler maps. For each scheduler map, the command output includes:

- The name of the scheduler map (scheduler-map field)
- The name of the scheduler (scheduler field)
- The forwarding classes mapped to the scheduler (forwarding-class field)
- The minimum guaranteed queue bandwidth (transmit-rate field)
- The scheduling priority (priority field)
- The maximum bandwidth in the priority group the queue can consume (shaping-rate field)
- The drop profile loss priority (loss priority field) for each drop profile name (name field)

The command output shows that:

- The scheduler map `be-map` was created and has these properties:
  - There are two schedulers, `be-sched` and `nc-sched`.
  - The scheduler `be-sched` has two forwarding classes, `best-effort` and `be2`.
  - Scheduler `be-sched` forwarding classes `best-effort` and `be2` share a minimum guaranteed bandwidth of 3,000,000,000 bps, can consume a maximum of 100 percent of the priority group bandwidth, and use the drop profile `dp-be-low` for low loss-priority traffic, the default drop profile for medium-high loss-priority traffic, and the drop profile `dp-be-high` for high loss-priority traffic.
  - The scheduler `nc-sched` has one forwarding class, `network-control`.
  - The `network-control` forwarding class has a minimum guaranteed bandwidth of 500,000,000 bps, can consume a maximum of 100 percent of the priority group bandwidth, and uses the drop profile `dp-nc` for low loss-priority traffic and the default drop profile for medium-high and high loss priority traffic.
- The scheduler map `gd-map` was created and has these properties:
  - There are two schedulers, `fcoe-sched` and `n1-sched`.
  - The scheduler `fcoe-sched` has one forwarding class, `fcoe`.
  - The `fcoe` forwarding class has a minimum guaranteed bandwidth of 2,500,000,000 bps, and can consume a maximum of 100 percent of the priority group bandwidth.
  - The scheduler `n1-sched` has one forwarding class, `no-loss`.
  - The `no-loss` forwarding class has a minimum guaranteed bandwidth of 2,000,000,000 bps, and can consume a maximum of 100 percent of the priority group bandwidth.
- The scheduler map `hpc-map` was created and has these properties:
  - There is one scheduler, `hpc-sched`.
  - The scheduler `hpc-sched` has one forwarding class, `hpc`.
  - The `hpc` forwarding class has a minimum guaranteed bandwidth of 2,000,000,000 bps, can consume a maximum of 100 percent of the priority group bandwidth, and uses the drop profile `dp-hpc` for low loss-priority traffic and the default drop profile for medium-high and high loss-priority traffic.

## Verifying the Drop Profiles

### Purpose

Verify that you created the drop profiles dp-be-high, dp-be-low, dp-hpc, and dp-nc with the correct fill levels and drop probabilities.

### Action

List the drop profiles using the operational mode command `show configuration class-of-service drop-profiles`:

```
user@switch> show configuration class-of-service drop-profiles
dp-be-low {
 interpolate {
 fill-level [25 50];
 drop-probability [0 80];
 }
}
dp-be-high {
 interpolate {
 fill-level [10 40];
 drop-probability [0 100];
 }
}
dp-hpc {
 interpolate {
 fill-level [75 90];
 drop-probability [0 75];
 }
}
dp-nc {
 interpolate {
 fill-level [80 100];
 drop-probability [0 100];
 }
}
```



## Meaning

The `show configuration class-of-service drop-profiles` command lists the drop profiles and their properties. The command output shows that there are four drop profiles configured, `dp-be-high`, `dp-be-low`, `dp-hpc`, and `dp-nc`. The output also shows that:

- For `dp-be-low`, the drop start point (the first fill level) is when the queue is 25 percent filled, the drop end point (the second fill level) occurs when the queue is 50 percent filled, and the drop probability at the drop end point is 80 percent.
- For `dp-be-high`, the drop start point (the first fill level) is when the queue is 10 percent filled, the drop end point (the second fill level) occurs when the queue is 40 percent filled, and the drop probability at the drop end point is 100 percent.
- For `dp-hpc`, the drop start point (the first fill level) is when the queue is 75 percent filled, the drop end point (the second fill level) occurs when the queue is 90 percent filled, and the drop probability at the drop end point is 75 percent.
- For `dp-nc`, the drop start point (the first fill level) is when the queue is 80 percent filled, the drop end point (the second fill level) occurs when the queue is 100 percent filled, and the drop probability at the drop end point is 100 percent.

## Verifying the Priority Group Output Schedulers (Traffic Control Profiles)

### Purpose

Verify that you created the traffic control profiles `be-tcp`, `gd-tcp`, and `hpc-tcp` with the correct bandwidth parameters and scheduler mapping.

### Action

List the traffic control profiles using the operational mode command `show class-of-service traffic-control-profile`:

```
user@switch> show class-of-service traffic-control-profile
Traffic control profile: be-tcp, Index: 40535
 Shaping rate: 100 percent
 Scheduler map: be-map
 Guaranteed rate: 3500000000

Traffic control profile: gd-tcp, Index: 37959
 Shaping rate: 100 percent
 Scheduler map: gd-map
```

```

Guaranteed rate: 4500000000

Traffic control profile: hpc-tcp, Index: 47661
 Shaping rate: 100 percent
 Scheduler map: hpc-map
 Guaranteed rate: 2000000000

```

## Meaning

The `show class-of-service traffic-control-profile` command lists all of the configured traffic control profiles. For each traffic control profile, the command output includes:

- The name of the traffic control profile (traffic-control-profile)
- The maximum port bandwidth the priority group can consume (shaping-rate)
- The scheduler map associated with the traffic control profile (scheduler-map)
- The minimum guaranteed priority group port bandwidth (guaranteed-rate)

The command output shows that:

- The traffic control profile `be-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `be-map`, and has a minimum guaranteed bandwidth of 3,500,000,000 bps.
- The traffic control profile `gd-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `gd-map`, and has a minimum guaranteed bandwidth of 4,500,000,000 bps.
- The traffic control profile `hpc-tcp` can consume a maximum of 100 percent of the port bandwidth, is associated with the scheduler map `hpc-map`, and has a minimum guaranteed bandwidth of 2,000,000,000 bps.

## Verifying the Interface Configuration

### Purpose

Verify that the classifier, the congestion notification profile, and the forwarding class sets are configured on interfaces `xe-0/0/20` and `xe-0/0/21`.

## Action

List the interfaces using the operational mode commands `show configuration class-of-service interfaces xe-0/0/20` and `show configuration class-of-service interfaces xe-0/0/21`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/20
forwarding-class-set {
 best-effort-gp {
 output-traffic-control-profile be-tcp;
 }
 guar-delivery-pg {
 output-traffic-control-profile gd-tcp;
 }
 hpc-pg {
 output-traffic-control-profile hpc-tcp;
 }
}
congestion-notification-profile gd_cnp;
unit 0 {
 classifiers {
 ieee-802.1 hsclassifier1;
 }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/21
forwarding-class-set {
 best-effort-gp {
 output-traffic-control-profile be-tcp;
 }
 guar-delivery-pg {
 output-traffic-control-profile gd-tcp;
 }
 hpc-pg {
 output-traffic-control-profile hpc-tcp;
 }
}
congestion-notification-profile gd_cnp;
unit 0 {
 classifiers {
 ieee-802.1 hsclassifier1;
 }
}
```

```
}
}
```

## Meaning

The `show configuration class-of-service interfaces interface-name` command shows that each interface includes the forwarding class sets `best-effort-pg`, `guar-delivery-pg`, and `hpc-pg`, congestion notification profile `gd-cnp`, and the IEEE 802.1p classifier `hsclassifier1`.

## RELATED DOCUMENTATION

*Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)*

[Benefits of Configuring CoS Hierarchical Port Scheduling](#)

*Assigning CoS Components to Interfaces*

*Example: Configuring WRED Drop Profiles*

*Example: Configuring Drop Profile Maps*

*Example: Configuring Forwarding Classes*

*Example: Configuring Forwarding Class Sets*

*Example: Configuring Queue Schedulers*

*Example: Configuring Queue Scheduling Priority*

*Example: Configuring Traffic Control Profiles (Priority Group Scheduling)*

*Example: Configuring Minimum Guaranteed Output Bandwidth*

*Example: Configuring Maximum Output Bandwidth*

*Configuring CoS PFC (Congestion Notification Profiles)*

[Overview of CoS Changes Introduced in Junos OS Release 12.2 | 67](#)

*Understanding CoS Hierarchical Port Scheduling (ETS)*

*Understanding CoS Scheduling Behavior and Configuration Considerations*

[Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports | 174](#)

[Understanding Default CoS Scheduling on QFabric System Interconnect Devices \(Junos OS Release 13.1 and Later Releases\) | 315](#)

## Disabling the ETS Recommendation TLV

The enhanced transmission selection (ETS) Recommendation TLV communicates the ETS settings that the switch wants the connected peer interface to use. If the peer interface is “willing,” the peer interface changes its configuration to match the configuration in the ETS Recommendation TLV. By default, the switch interfaces send the ETS Recommendation TLV to the peer. The settings communicated are the egress ETS settings defined by configuring hierarchical scheduling on the interface.

We recommend that you use the same ETS settings on the connected peer that you use on the switch interface and that you leave the ETS Recommendation TLV enabled. However, on interfaces that use IEEE DCBX as the DCBX mode, if you want an asymmetric configuration between the switch interface and the connected peer, you can disable the ETS Recommendation TLV.

**NOTE:** Disabling the ETS Recommendation TLV on interfaces that use DCBX version 1.01 as the DCBX mode has no effect and does not change DCBX behavior.

If you disable the ETS Recommendation TLV, the switch still sends the ETS Configuration TLV to the connected peer. The result is that the connected peer is informed about the switch DCBX ETS configuration, but even if the peer is “willing,” the peer does not change its configuration to match the switch configuration. This is asymmetric configuration—the two interfaces can have different parameter values for the ETS attribute.

To disable the ETS Recommendation TLV:

- ```
[edit protocols dcbx interface interface-name]  
user@switch# set enhanced-transmission-selection no-recommendation-tlv
```

RELATED DOCUMENTATION

Configuring the DCBX Mode

Configuring DCBX Autonegotiation

Understanding DCBX

[Understanding Data Center Bridging Capability Exchange Protocol for EX Series Switches](#)

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

IN THIS SECTION

- [General Information about Ethernet PAUSE and PFC and When to Use Them | 492](#)
- [Ethernet PAUSE | 492](#)
- [PFC | 498](#)
- [Lossless Transport Support Summary | 502](#)

Flow control supports lossless transmission by regulating traffic flows to avoid dropping frames during periods of congestion. Flow control stops and resumes the transmission of network traffic between two connected peer nodes on a full-duplex Ethernet physical link. Controlling the flow by pausing and restarting it prevents buffers on the nodes from overflowing and dropping frames. You configure flow control on a per-interface basis.

Two methods of peer-to-peer flow control are supported:

- IEEE 802.3X Ethernet PAUSE

NOTE: QFX10000 switches do not support Ethernet PAUSE. Information about Ethernet PAUSE does not apply to QFX10000 switches.

OCX Series switches support symmetric Ethernet PAUSE flow control on Layer 3 tagged interfaces. OCX Series switches do not support asymmetric Ethernet PAUSE flow control. Information about asymmetric flow control does not apply to OCX Series switches.

- IEEE 802.1Qbb priority-based flow control (PFC)

NOTE: OCX Series switches do not support PFC or lossless Layer 2 transport. Information about PFC, lossless transport, and congestion notification profiles does not apply to OCX Series switches.

NOTE: QFX10002-60C devices do not support PFC and lossless queues; that is, the default lossless queues (fcoe and no-loss) will be lossy queues.



Video: [Why Use PFC in a Data Center Network?](#)

General Information about Ethernet PAUSE and PFC and When to Use Them

Ethernet PAUSE and PFC are link-level flow control mechanisms.

NOTE: For end-to-end congestion control for best-effort traffic, see [Understanding CoS Explicit Congestion Notification](#).

Ethernet PAUSE pauses transmission of all traffic on a physical Ethernet link.

PFC decouples the pause function from the physical Ethernet link and enables you to divide traffic on one link into eight priorities. You can think of the eight priorities as eight “lanes” of traffic that are mapped to forwarding classes and output queues. Each priority maps to a 3-bit IEEE 802.1p CoS code point value in the VLAN header. You can enable PFC on one or more priorities (IEEE 802.1p code points) on a link. When PFC-enabled traffic is paused on a link, traffic that is not PFC-enabled continues to flow (or is dropped if congestion is severe enough).

Use Ethernet PAUSE when you want to prevent packet loss on all of the traffic on a link. Use PFC to prevent traffic loss only on a specified type of traffic that require lossless treatment, for example, Fibre Channel over Ethernet (FCoE) traffic.

NOTE: Depending on the amount of traffic on a link or assigned to a priority, pausing traffic can cause ingress port congestion and spread congestion through the network.

Ethernet PAUSE and PFC are mutually exclusive configurations on an interface. Attempting to configure both Ethernet PAUSE and PFC on a link causes a commit error.

By default, all forms of flow control are disabled. You must explicitly enable flow control on interfaces to pause traffic.

Ethernet PAUSE

Ethernet PAUSE is a congestion relief feature that works by providing link-level flow control for all traffic on a full-duplex Ethernet link. Ethernet PAUSE works in both directions on the link. In one direction, an interface generates and sends Ethernet PAUSE messages to stop the connected peer from sending more traffic. In the other direction, the interface responds to Ethernet PAUSE messages it receives from the connected peer to stop sending traffic.

NOTE: QFX10000 switches do not support Ethernet PAUSE. Information about Ethernet PAUSE does not apply to QFX10000 switches.

OCX Series switches support symmetric Ethernet PAUSE flow control on Layer 3 tagged interfaces. OCX Series switches do not support asymmetric Ethernet PAUSE flow control. Information about asymmetric flow control does not apply to OCX Series switches.

Ethernet PAUSE also works on aggregated Ethernet interfaces. For example, if the connected peer interfaces are called Node A and Node B:

- When the receive buffers on interface Node A reach a certain level of fullness, the interface generates and sends an Ethernet PAUSE message to the connected peer (interface Node B) to tell the peer to stop sending frames. The Node B buffers store frames until the time period specified in the Ethernet PAUSE frame elapses; then Node B resumes sending frames to Node A.
- When interface Node A receives an Ethernet PAUSE message from interface Node B, interface Node A stops transmitting frames until the time period specified in the Ethernet PAUSE frame elapses; then Node A resumes transmission. (The Node A transmit buffers store frames until Node A resumes sending frames to Node B.)

In this scenario, if Node B sends an Ethernet PAUSE frame with a time value of 0 to Node A, the 0 time value indicates to Node A that it can resume transmission. This happens when the Node B buffer empties to below a certain threshold and the buffer can once again accept traffic.

Symmetric flow control means an interface has the same Ethernet PAUSE configuration in both directions. The Ethernet PAUSE generation and Ethernet PAUSE response functions are both configured as enabled, or they are both disabled. You configure symmetric flow control by including the `flow-control` statement at the `[edit interfaces interface-name ether-options]` hierarchy level.

Asymmetric flow control allows you to configure the Ethernet PAUSE functionality in each direction independently on an interface. The configuration for generating Ethernet PAUSE messages and for responding to Ethernet PAUSE messages does not have to be the same. It can be enabled in both directions, disabled in both directions, or enabled in one direction and disabled in the other direction. You configure asymmetric flow control by including the `configured-flow-control` statement at the `[edit interfaces interface-name ether-options]` hierarchy level.

On any particular interface, symmetric and asymmetric flow control are mutually exclusive. Asymmetric flow control overrides and disables symmetric flow control. (If PFC is configured on an interface, you cannot commit an Ethernet PAUSE configuration on the interface. Attempting to commit an Ethernet PAUSE configuration on an interface with PFC enabled on one or more queues results in a commit error. To commit the PAUSE configuration, you must first delete the PFC configuration.) Both symmetric and asymmetric flow control are supported.

Symmetric Flow Control

Symmetric flow control configures both the receive and transmit buffers in the same state. The interface can both send Ethernet PAUSE messages and respond to them (flow control is enabled), or the interface cannot send Ethernet PAUSE messages or respond to them (flow control is disabled).

When you enable symmetric flow control on an interface, the Ethernet PAUSE behavior depends on the configuration of the connected peer. With symmetric flow control enabled, the interface can perform any Ethernet PAUSE functions that the connected peer can perform. (When symmetric flow control is disabled, the interface does not send or respond to Ethernet PAUSE messages.)

Asymmetric Flow Control

Asymmetric flow control enables you to specify independently whether or not the interface receive buffer generates and sends Ethernet PAUSE messages to stop the connected peer from transmitting traffic, and whether or not the interface transmit buffer responds to Ethernet PAUSE messages it receives from the connected peer and stops transmitting traffic. The receive buffer configuration determines if the interface transmits Ethernet PAUSE messages, and the transmit buffer configuration determines if the interface receives and responds to Ethernet PAUSE messages:

- Receive buffers on—Enable Ethernet PAUSE transmission (generate and send Ethernet PAUSE frames)
- Transmit buffers on—Enable Ethernet PAUSE reception (respond to received Ethernet PAUSE frames)

You must explicitly set the flow control for both the receive buffer and the transmit buffer (on or off) to configure asymmetric Ethernet PAUSE. [Table 86 on page 494](#) describes the configured flow control state when you set the receive (Rx) and transmit (Tx) buffers on an interface:

Table 86: Asymmetric Ethernet PAUSE Flow Control Configuration

Receive (Rx) Buffer	Transmit (Tx) Buffer	Configured Flow Control State
On	Off	Interface generates and sends Ethernet PAUSE messages. Interface does not respond to Ethernet PAUSE messages (interface continues to transmit even if peer requests that the interface stop sending traffic).
Off	On	Interface responds to Ethernet PAUSE messages received from the connected peer, but does not generate or send Ethernet PAUSE messages. (The interface does not request that the connected peer stop sending traffic.)

Table 86: Asymmetric Ethernet PAUSE Flow Control Configuration (Continued)

Receive (Rx) Buffer	Transmit (Tx) Buffer	Configured Flow Control State
On	On	Same functionality as symmetric Ethernet PAUSE. Interface generates and sends Ethernet PAUSE messages and responds to received Ethernet PAUSE messages.
Off	Off	Ethernet PAUSE flow control is disabled.

The configured flow control is the Ethernet PAUSE state configured on the interface.

On 1-Gigabit Ethernet interfaces, autonegotiation of Ethernet PAUSE with the connected peer is supported. (Autonegotiation on 10-Gigabit Ethernet interfaces is not supported.) Autonegotiation enables the interface to exchange state advertisements with the connected peer so that the two devices can agree on the Ethernet PAUSE configuration. Each interface advertises its flow control state to the connected peer using a combination of the Ethernet PAUSE and ASM_DIR bits, as described in [Table 87 on page 495](#):

Table 87: Flow Control State Advertised to the Connected Peer (Autonegotiation)

Rx Buffer State	Tx Buffer State	PAUSE Bit	ASM_DIR Bit	Description
Off	Off	0	0	The interface advertises no Ethernet PAUSE capability. This is equivalent to disabling flow control on an interface.
On	On	1	0	The interface advertises symmetric flow control (both the transmission of Ethernet PAUSE messages and the ability to receive and respond to Ethernet PAUSE messages).

Table 87: Flow Control State Advertised to the Connected Peer (Autonegotiation) (*Continued*)

Rx Buffer State	Tx Buffer State	PAUSE Bit	ASM_DIR Bit	Description
On	Off	0	1	The interface advertises asymmetric flow control (the transmission of Ethernet PAUSE messages, but not the ability to receive and respond to Ethernet PAUSE messages).
Off	On	1	1	The interface advertises both symmetric and asymmetric flow control. Although the interface does not generate and send Ethernet PAUSE requests to the peer, the interface supports both symmetric and asymmetric Ethernet PAUSE configuration on the peer because the peer is not affected if the peer does not receive Ethernet PAUSE requests. (If the interface responds to the peer's Ethernet PAUSE requests, that is sufficient to support either symmetric or asymmetric flow control on the peer.)

The flow control configuration on each switch interface interacts with the flow control configuration of the connected peer. Each peer advertises its state to the other peer. The interaction of the flow control configuration of the peers determines the flow control behavior (resolution) between them, as shown in [Table 88 on page 497](#). The first four columns show the Ethernet PAUSE configuration on the local QFX Series or EX4600 switch and on the connected peer (also known as the *link partner*). The last two columns show the Ethernet PAUSE resolution that results from the local and peer configurations on each interface. This illustrates how the Ethernet PAUSE configuration of each interface affects the Ethernet PAUSE behavior on the other interface.

NOTE: In the Resolution columns of the table, disabling Ethernet PAUSE transmit means that the interface receive buffers do not generate and send Ethernet PAUSE messages to the peer. Disabling Ethernet PAUSE receive means that the interface transmit buffers do not respond to Ethernet PAUSE messages received from the peer.

Table 88: Asymmetric Ethernet PAUSE Behavior on Local and Peer Interfaces

Local Interface (QFX Series or EX4600 Switch)		Peer Interface		Local Resolution	Peer Resolution
PAUSE Bit	ASM_DIR Bit	PAUSE Bit	ASM_DIR Bit		
0	0	Don't care	Don't care	Disable Ethernet PAUSE transmit and receive	Disable Ethernet PAUSE transmit and receive
0	1	0	Don't care	Disable Ethernet PAUSE transmit and receive	Disable Ethernet PAUSE transmit and receive
0	1	1	0	Disable Ethernet PAUSE transmit and receive	Disable Ethernet PAUSE transmit and receive
0	1	1	1	Enable Ethernet PAUSE transmit and disable Ethernet PAUSE receive	Disable Ethernet PAUSE transmit and enable Ethernet PAUSE receive
1	0	0	Don't care	Disable Ethernet PAUSE transmit and receive	Disable Ethernet PAUSE transmit and receive
1	0	1	Don't care	Enable Ethernet PAUSE transmit and receive	Enable Ethernet PAUSE transmit and receive
1	1	0	0	Disable Ethernet PAUSE transmit and receive	Disable Ethernet PAUSE transmit and receive

Table 88: Asymmetric Ethernet PAUSE Behavior on Local and Peer Interfaces (Continued)

Local Interface (QFX Series or EX4600 Switch)		Peer Interface		Local Resolution	Peer Resolution
PAUSE Bit	ASM_DIR Bit	PAUSE Bit	ASM_DIR Bit		
1	1	0	1	Enable Ethernet PAUSE receive and disable Ethernet PAUSE transmit	Enable Ethernet PAUSE transmit and disable Ethernet PAUSE receive
1	1	Don't care	Don't care	Enable Ethernet PAUSE transmit and receive	Enable Ethernet PAUSE transmit and receive

NOTE: For your convenience, [Table 88 on page 497](#) replicates Table 28B-3 of Section 2 of the IEEE 802.X specification.

PFC

PFC is a lossless transport and congestion relief feature that works by providing granular link-level flow control for each IEEE 802.1p code point (priority) on a full-duplex Ethernet link. When the receive buffer on a switch interface fills to a threshold, the switch transmits a pause frame to the sender (the connected peer) to temporarily stop the sender from transmitting more frames. The buffer threshold must be low enough so that the sender has time to stop transmitting frames and the receiver can accept the frames already on the wire before the buffer overflows. The switch automatically sets queue buffer thresholds to prevent frame loss.

When congestion forces one priority on a link to pause, all of the other priorities on the link continue to send frames. Only frames of the paused priority are not transmitted. When the receive buffer empties below another threshold, the switch sends a message that starts the flow again.

You configure PFC using a congestion notification profile (CNP). A CNP has two parts:

- **Input**—Specify the code point (or code points) on which to enable PFC, and optionally specify the maximum receive unit (MRU) and the cable length between the interface and the connected peer interface.
- **Output**—Specify the output queue or output queues that respond to pause messages from the connected peer.

You apply a PFC configuration by configuring a CNP on one or more interfaces. Each interface that uses a particular CNP is enabled to pause traffic identified by the priorities (code points) specified in that CNP. You can configure one CNP on an interface, and you can configure different CNPs on different interfaces. When you configure a CNP on an interface, ingress traffic that is mapped to a priority that the CNP enables for PFC is paused whenever the queue buffer fills to the pause threshold. (The pause threshold is not user-configurable.)

Configure PFC for a priority end to end along the entire data path to create a lossless lane of traffic on the network. You can selectively pause the traffic in any queue without pausing the traffic for other queues on the same link. You can create lossless lanes for traffic such as FCoE, LAN backup, or management, while using standard frame-drop congestion management for IP traffic on the same link.

Potential consequences of flow control are:

- Ingress port congestion (configuring too many lossless flows can cause ingress port congestion)
- A paused priority that causes upstream devices to pause the same priority, thus spreading congestion back through the network

By definition, PFC supports symmetric pause only (as opposed to Ethernet PAUSE, which supports symmetric and asymmetric pause). With symmetric pause, a device can:

- Transmit pause frames to pause incoming traffic. (You configure this using the input stanza of a congestion notification profile.)
- Receive pause frames and stop sending traffic to a device whose buffer is too full to accept more frames. (You configure this using the output stanza of a congestion notification profile.)

Receiving a PFC frame from a connected peer pauses traffic on egress queues based on the IEEE 802.1p priorities that the PFC pause frame identifies. The priorities are 0 through 7. By default, the priorities map to queue numbers 0 through 7, respectively, and to specific forwarding classes, as shown in [Table 89 on page 499](#):

Table 89: Default PFC Priority to Queue and Forwarding Class Mapping

IEEE 802.1p Priority (Code Point)	Queue	Forwarding Class
0 (000)	0	best-effort
1 (001)	1	best-effort
2 (010)	2	best-effort

Table 89: Default PFC Priority to Queue and Forwarding Class Mapping (Continued)

IEEE 802.1p Priority (Code Point)	Queue	Forwarding Class
3 (011)	3	fcoe
4 (100)	4	no-loss
5 (101)	5	best-effort
6 (110)	6	network-control
7 (111)	7	network-control

For example, a received PFC pause frame that pauses priority 3 pauses output queue 3. If you do not want to use the default configuration, you can configure customized mapping of priorities to queues and forwarding classes.

NOTE: By convention, deployments with converged server access typically use IEEE 802.1p priority 3 for FCoE traffic. The default configuration sets the `fcoe` forwarding class as a lossless forwarding class that is mapped to queue 3. The default classifier maps incoming priority 3 traffic to the `fcoe` forwarding class. *However, you must apply PFC to the entire FCoE data path to configure the end-to-end lossless behavior that FCoE traffic requires.*

If your network uses priority 3 for FCoE traffic, we recommend that you use the default configuration. If your network uses a priority other than 3 for FCoE traffic, you can configure lossless FCoE transport on any IEEE 802.1p priority as described in [Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows](#) and [Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway](#).

To enable PFC on a priority:

1. Specify the IEEE 802.1p code point to pause in the input stanza of a CNP.
2. If you are not using the default lossless forwarding classes, specify the IEEE 802.1p code point to pause and the corresponding output queue in the output stanza of the CNP.
3. Apply the CNP to the ingress interfaces on which you want to pause the traffic.
4. If you are not using the default lossless forwarding classes, apply the CNP to the ingress interfaces on which you want to pause the traffic.



CAUTION: Any change to the PFC configuration on a port temporarily blocks the entire port (not just the priorities affected by the PFC change) so that the port can implement the change, then unblocks the port. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

A change to the PFC configuration means any change to a CNP, including changing the input portion of the CNP (enabling or disabling PFC on a priority, or changing the MRU or cable-length values) or changing the output portion of the CNP that enables or disables output flow control on a queue. A PFC configuration change only affects ports that use the changed CNP.

The following actions change the PFC configuration:

- Deleting or disabling a PFC configuration (input or output) in a CNP that is in use on one or more interfaces. For example:
 1. An existing CNP with an input stanza that enables PFC on priorities 3, 5, and 6 is configured on interfaces xe-0/0/20 and xe-0/0/21.
 2. We disable the PFC configuration for priority 6 in the input CNP, and then commit the configuration.
 3. The PFC configuration change causes all traffic on interfaces xe-0/0/20 and xe-0/0/21 to stop until the PFC change has been implemented. When the PFC change has been implemented, traffic resumes.
- Configuring a CNP on an interface. (This changes the PFC state by enabling PFC on one or more priorities.)
- Deleting a CNP from an interface. (This changes the PFC state by disabling PFC on one or more priorities.)

When you associate the CNP with an interface, the interface uses PFC to send pause requests when the output queue buffer for the lossless traffic fills to the pause threshold.

On switches that use different classifiers for unicast and multidestination traffic, you can map a unicast queue (queue 0 through 7) and a multidestination queue (queue 8, 9, 10, or 11) to the same IEEE 802.1p code point (priority) so that both unicast and multicast traffic use that priority. However, do not map multidestination traffic to lossless output queues. Starting with Junos OS Release 12.3, you can map one priority to multiple output queues.

NOTE: You can attach a maximum of one CNP to an interface, but you can create an unlimited number of CNPs that explicitly configure only the input stanza and use the default output stanza.

The output stanza of the CNP maps to a profile that interfaces use to respond to pause messages received from the connected peer. On standalone switches, you can create two CNPs with an explicitly configured output stanza.

When a switch is a Node device in a QFabric system, you can create one CNP with an explicitly configured output stanza. (One fewer profile is available on QFabric systems because the system needs a default profile for fabric interfaces, which are not used as fabric interfaces when the switches are not part of a QFabric system. [Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows](#) describes configuring output flow control.

Lossless Transport Support Summary

The switch supports up to six lossless forwarding classes. For lossless transport, you must enable PFC on the IEEE 802.1p priorities (code points) mapped to lossless forwarding classes.



CAUTION: Any change to the PFC configuration on a port temporarily blocks the entire port (not just the priorities affected by the PFC change) so that the port can implement the change, then unblocks the port. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

The following limitation applies to support lossless transport on QFabric systems only:

- The internal fiber cable length from the QFabric system Node device to the QFabric system Interconnect device cannot exceed 150 meters.

The default CoS configuration provides two lossless forwarding classes, *fcoe* and *no-loss*. If you explicitly configure lossless forwarding classes, you must include the *no-loss* packet drop attribute to enable lossless behavior, or the traffic is not lossless. For both default and explicit lossless forwarding class configuration, you must configure CNP input stanzas to enable PFC on the priority of the lossless traffic and apply the CNPs to ingress interfaces.

NOTE: The information in this note applies only to systems that do not run the ELS CLI.

Junos OS Release 12.2 introduced changes to the way the switch handles lossless forwarding classes (including the default *fcoe* and *no-loss* forwarding classes).

In Junos OS Release 12.1, either explicitly configuring the *fcoe* and *no-loss* forwarding classes or using the default configuration for these forwarding classes resulted in the same lossless behavior for traffic mapped to those forwarding classes.

However, in Junos OS Release 12.2, if you explicitly configure the `fcoe` or the `no-loss` forwarding class, that forwarding class is no longer treated as a lossless forwarding class. Traffic mapped to these forwarding classes is treated as lossy (best-effort) traffic. This is true even if the explicit configuration is exactly the same as the default configuration.

If your CoS configuration from Junos OS Release 12.1 or earlier includes the explicit configuration of the `fcoe` or the `no-loss` forwarding class, then when you upgrade to Junos OS Release 12.2, those forwarding classes are not lossless. To preserve the lossless treatment of these forwarding classes, delete the the explicit `fcoe` and `no-loss` forwarding class configuration before you upgrade to Junos OS Release 12.2.

See [Overview of CoS Changes Introduced in Junos OS Release 12.2](#) for detailed information about this change and how to delete an existing lossless configuration.

In Junos OS Release 12.3, the default behavior of the `fcoe` and `no-loss` forwarding classes is the same as in Junos OS Release 12.2. However, in Junos OS Release 12.3, you can configure up to six lossless forwarding classes. All explicitly configured lossless forwarding classes must include the new `no-loss` packet drop attribute or the forwarding class is lossy.

[Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows](#) provides detailed information about the explicit configuration of lossless priorities and about the default configuration of lossless priorities, including the input and output stanzas of the CNP.

NOTE: PFC and Ethernet PAUSE are used only on Ethernet interfaces. Fabric (fte) ports on QFabric systems (Node device fabric ports and Interconnect device fabric ports) use link-layer flow control (LLFC) to ensure the appropriate treatment of lossless traffic.

Release History Table

Release	Description
21.2R1EVO	PTX10008 routers support DCBX and PFC.
12.3	Starting with Junos OS Release 12.3, you can map one priority to multiple output queues.

RELATED DOCUMENTATION

Understanding DCB Features and Requirements

[Understanding CoS Explicit Congestion Notification](#)

[Configuring CoS PFC \(Congestion Notification Profiles\)](#)

Enabling and Disabling CoS Symmetric Ethernet PAUSE Flow Control

Ethernet PAUSE flow control is a congestion relief feature that works by providing link-level flow control for all traffic on a full-duplex Ethernet link, including Ethernet links that belong to Ethernet link aggregated (LAG) interfaces. Ethernet PAUSE works in both directions on the link. In one direction, an interface generates and sends PAUSE messages to stop the connected peer from sending more traffic. In the other direction, the interface responds to PAUSE messages it receives from the connected peer to stop sending traffic.

Symmetric flow control means that an interface has the same PAUSE configuration in both directions. The PAUSE generation and PAUSE response functions are both configured as enabled, or they are both disabled.

Asymmetric flow control allows you to configure the PAUSE functionality in each direction independently on an interface. The configuration for generating PAUSE messages and for responding to PAUSE messages does not have to be the same. It can be enabled in both directions, disabled in both directions, or enabled in one direction and disabled in the other direction. If you do not want to PAUSE all of the traffic on a link, you can use priority-based flow control (PFC) to selectively pause traffic based on its IEEE 802.1p code point.

NOTE: OCX Series switches do not support PFC.

On any particular interface, symmetric and asymmetric flow control are mutually exclusive. If you attempt to configure both features, the switch returns a commit error. Ethernet PAUSE and PFC are also mutually exclusive features, so you cannot configure both of them on the same interface. If you attempt to configure both Ethernet PAUSE and PFC on an interface, the switch returns a commit error.

By default, all flow control features are disabled. You enable symmetric flow control on the interfaces on which you want to PAUSE all of the traffic on a link.

- To enable symmetric flow control on an interface:

```
[edit interfaces interface-name ether-options]  
user@switch# set flow-control
```

- To disable symmetric flow control on an interface:

```
[edit interfaces interface-name ether-options]  
user@switch# set no-flow-control
```

RELATED DOCUMENTATION

Configuring CoS Asymmetric Ethernet PAUSE Flow Control

Configuring CoS PFC (Congestion Notification Profiles)

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Configuring CoS Asymmetric Ethernet PAUSE Flow Control

Ethernet PAUSE flow control is a congestion relief feature that works by providing link-level flow control for all traffic on a full-duplex Ethernet link, including Ethernet links that belong to link aggregated (LAG) interfaces. Ethernet PAUSE works in both directions on the link. In one direction, an interface generates and sends PAUSE messages to stop the connected peer from sending more traffic. In the other direction, the interface responds to PAUSE messages it receives from the connected peer to stop sending traffic.

Asymmetric flow control allows you to configure the PAUSE functionality in each direction independently on an interface. The configuration for generating PAUSE messages and for responding to PAUSE messages does not have to be the same. It can be enabled in both directions, disabled in both directions, or enabled in one direction and disabled in the other direction.

Symmetric flow control means that the interface has the same configuration in both directions. The PAUSE generation and PAUSE response functions are both configured as enabled or they are both disabled. If you do not want to PAUSE all of the traffic on a link, you can use priority-based flow control (PFC) to selectively pause traffic based on its IEEE 802.1p code point.

Asymmetric flow control provides the ability to configure the receive buffer and transmit buffer Ethernet PAUSE actions independently on an interface. The buffers perform the following actions:

- The receive buffers generate and send PAUSE messages to the connected peer to ask the peer to stop sending traffic for a time period specified in the PAUSE frame. The peer interface's buffers may store outgoing frames until the PAUSE period elapses and the interface can resume sending traffic.
- The transmit buffers respond to PAUSE messages received from the connected peer to stop sending traffic to the peer. The transmit buffer may store outgoing frames until the PAUSE period elapses and the interface can resume sending traffic.

Asymmetric flow control enables you to specify independently whether or not the interface receive buffer generates and sends PAUSE messages to stop the connected peer from transmitting traffic, and whether or not the interface transmit buffer responds to PAUSE messages it receives from the connected peer and stops transmitting traffic. The receive buffer configuration determines if the interface transmits PAUSE messages, and the transmit buffer configuration determines if the interface receives and responds to PAUSE messages:

- Receive buffers on—Enable PAUSE transmission (generate and send PAUSE frames)
- Transmit buffers on—Enable PAUSE reception (respond to received PAUSE frames)

You must explicitly set both the receive buffer and the transmit buffer to configure asymmetric flow control.

- To configure asymmetric flow control on an interface:

```
[edit interfaces interface-name ether-options]
user@switch# set configured-flow-control rx-buffers (on | off) tx-buffers (on | off)
```

For example, to configure interface xe-0/0/24 to generate and send PAUSE messages but not to respond to received PAUSE messages:

```
set interfaces xe-0/0/24 ether-options configured-flow-control rx-buffers on tx-buffers off
```

For example, to configure interface xe-0/0/30 to respond to received PAUSE messages but not to generate and send PAUSE messages:

```
set interfaces xe-0/0/30 ether-options configured-flow-control rx-buffers off tx-buffers on
```

NOTE: If you configure both buffers to be on, that is equivalent to symmetric flow control. If you configure both buffers to be off, there is no flow control (flow control is disabled).

RELATED DOCUMENTATION

Enabling and Disabling CoS Symmetric Ethernet PAUSE Flow Control

Configuring CoS PFC (Congestion Notification Profiles)

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Configuring CoS PFC (Congestion Notification Profiles)

A congestion notification profile (CNP) enables priority-based flow control (PFC) on specified IEEE 802.1p priorities (code points). A CNP has two components:

- Input CNP:
 - Enable PFC on a specified priority.
 - Configure the maximum receive unit (MRU) on an interface for traffic that matches the PFC priority (optional).
 - Specify the length of the attached cable on the ingress interface (optional)
- Output CNP (optional): Configure flow control to enable PFC pause on specific output queues for specified priorities.

NOTE: By default, output queues 3 and 4 (which are mapped to default lossless forwarding classes `fcoe` and `no-loss`, respectively) are configured to respond to PFC pause messages received from the connected peer on priorities 3 and 4 (code points 011 and 100, respectively). If you explicitly configure flow control on any output queue, you must configure flow control on every output queue that you want to respond to pause messages. (The explicit configuration overrides the default configuration.)

To achieve lossless behavior, the output queue priorities on which you enable PFC flow control must match the PFC priorities on which you enable PFC on the input interfaces. For example, if you program output queues to pause priorities 3 (011) and 5 (101) in the output component of the CNP, then you must also enable pause on priorities 3 and 5 on the input component of the CNP. (In addition, the forwarding classes mapped to the paused output queues must be lossless forwarding classes.)

Associating a CNP with an interface enables PFC on the ingress traffic that matches the priority specified in the input CNP, and programs the queues listed in the output CNP to pause when the interface receives a PFC pause message from the connected peer. Configure PFC on a priority end to end along the entire data path to create a lossless lane of traffic on the network.

NOTE: You must enable PFC on the priority used by FCoE traffic on ingress interfaces (input CNP). Enable PFC on the FCoE priority on every interface that carries FCoE traffic. By convention, FCoE traffic uses priority 3 (code point 011), which maps to queue 3. If your network uses priority 3 for FCoE traffic, the default forwarding class and classifier configuration support

lossless transport, but you must still configure a CNP and apply it to the correct ingress interfaces to enable PFC and achieve lossless transport.

If your network does not use priority 3 for FCoE traffic, you need to configure a classifier that classifies FCoE traffic into a lossless forwarding class, based on the priority your network uses for FCoE traffic. If you are not using the default lossless forwarding class configuration, then you also need to ensure that the output queue mapped to the lossless FCoE forwarding class is programmed to pause.

You can attach only one CNP to an interface. There is no limit to the total number of CNPs you can create.

Configuring a CNP consists of:

- Naming the CNP.
- Specifying the IEEE 802.1 code point (priority) on which you want to enable PFC on ingress interfaces (input CNP).
- Optionally, specifying the MRU and the length of the attached cable on ingress interfaces (input CNP).
- Optionally, configuring flow control (PFC pause) on specified output queues if you want queues other than queues 3 and 4 to respond to pause messages received from the connected peer (output CNP).
- Mapping the CNP to an interface.

NOTE: Configuring or changing PFC on an interface blocks the entire port until the PFC change is completed. After a PFC change is completed, the port is unblocked and traffic resumes. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

NOTE: On QFX5100, QFX5200, and QFX5210, once the headroom buffer is exhausted, any new CNP configuration is not allocated headroom buffer, even if headroom buffer is freed by deletion of an existing CNP. CNP configuration has to be applied again to re-allocate the headroom buffer.



CAUTION: On QFX5130 and QFX5220, you must map all PFC-enabled IEEE 802.1P code-points to a lossless (no-loss) forwarding class. If a CNP has code-points that are

mapped to a lossy forwarding class, the entire CNP will not be programmed in hardware.

1. Enable PFC on the desired priority in the input CNP and optionally configure the interface MRU for traffic on that priority:

```
[edit class-of-service]
user@switch# set congestion-notification-profile cnp-name input ieee-802.1 code-point code-point bits pfc mru mru-value
```

For example, to configure a CNP named `fcoe-cnp` that enables PFC on IEEE 802.1 code point 011 and configures an MRU value of 2240:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
mru 2240
```

2. (Optional) Configure the length of the cable attached to the ingress interface:

```
[edit class-of-service]
user@switch# set congestion-notification-profile cnp-name input cable-length cable-length-value
```

For example, to configure a CNP named `fcoe-cnp` that sets the length of the ingress interface cable to 100 meters:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe-cnp input cable-length 100
```

3. (Optional) Configure flow control on output queues:

```
[edit class-of-service]
user@switch# set congestion-notification-profile cnp-name output ieee-802.1 code-point code-point-bits flow-control-queue [queue | list-of-queues]
```


For example, to configure a CNP named `fcoe-cnp` that enables PFC pause flow control on output queues 3 and 5 for FCoE traffic that uses priority 3 (code point 011) and on output queue 4 for traffic that uses priority 4 (code point 100):

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe-cnp output ieee-802.1 code-point 011
flow-control-queue [3 5]
user@switch# set congestion-notification-profile fcoe-cnp output ieee-802.1 code-point 100
flow-control-queue 4
```

4. Map the CNP to an interface:

```
[edit class-of-service]
user@switch# set interfaces interface congestion-notification-profile cnp-name
```

For example, to map the CNP `fcoe-cnp` to the interface `xe-0/0/7`:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/7 congestion-notification-profile fcoe-cnp
```

RELATED DOCUMENTATION

Example: Configuring CoS PFC for FCoE Traffic

Assigning CoS Components to Interfaces

Monitoring Interfaces That Have CoS Components

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows

Example: Configuring CoS PFC for FCoE Traffic

IN THIS SECTION

● [Requirements | 511](#)

- Overview | 511
- Configuration | 514
- Verification | 521

Priority-based flow control (PFC, described in IEEE 802.1Qbb) is a link-level flow control mechanism that you apply at ingress interfaces. PFC enables you to divide traffic on one physical link into eight priorities. You can think of the eight priorities as eight “lanes” of traffic that correspond to queues (forwarding classes). Each priority is mapped to a 3-bit IEEE 802.1p CoS value in the VLAN header.

You can selectively apply PFC to the traffic in any queue without pausing the traffic in other queues on the same link. You must apply PFC to FCoE traffic to ensure lossless transport.

This example describes how to configure PFC for FCoE traffic:

Requirements

This example uses the following hardware and software components:

- One switch
- Junos OS Release 11.1 or later for the QFX Series

Overview

IN THIS SECTION

- Topology | 512

FCoE traffic requires PFC to ensure lossless packet transport. This example shows you how to configure PFC on FCoE traffic, use the default FCoE forwarding-class-to-queue mapping and:

- Configure a classifier that associates the FCoE forwarding class with FCoE traffic, which is identified by IEEE 802.1p code point 011 (priority 3).
- Configure a congestion notification profile to apply PFC to the FCoE traffic.
- Apply the classifier and the PFC configuration to ingress interfaces.

NOTE: Configuring or changing PFC on an interface blocks the entire port until the PFC change is completed. After a PFC change is completed, the port is unblocked and traffic resumes. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

- Configure the CoS bandwidth scheduling for the FCoE forwarding class output queue.
- On switches that support enhanced transmission selection (ETS) hierarchical port scheduling, create a forwarding class set (priority group) that includes the FCoE forwarding class; this is required to configure enhanced transmission selection (ETS) and support data center bridging (DCB).
- For ETS, configure the bandwidth scheduling for the FCoE priority group.
- Apply the configuration to ingress and egress interfaces. How this is done differs depending on whether you use ETS or direct port scheduling for the CoS configuration.

For direct port scheduling, you apply a scheduler map directly to the interface. A scheduler map maps schedulers to forwarding classes, and applies the CoS properties of the scheduler to the output queue mapped to the forwarding class.

For ETS hierarchical port scheduling, you apply the scheduler map to a traffic control profile, and then apply the traffic control profile to the interface. The scheduler map maps CoS properties to forwarding classes (and their associated output queues) just as it does for direct port scheduling. The traffic control profile maps CoS properties to the priority group (a group of forwarding classes defined in a forwarding class set) that contains the forwarding class, creating a CoS hierarchy that allocates port bandwidth to a group of forwarding classes (priority group), and then allocates the priority group bandwidth to the individual forwarding classes.

Each interface in this example acts as both an ingress interface and an egress interface, so the classifier, congestion notification profile, and scheduling are applied to all of the interfaces.

Topology

Table 90 on page 512 shows the configuration components for this example.

Table 90: Components of the PFC for FCoE Traffic Configuration Topology

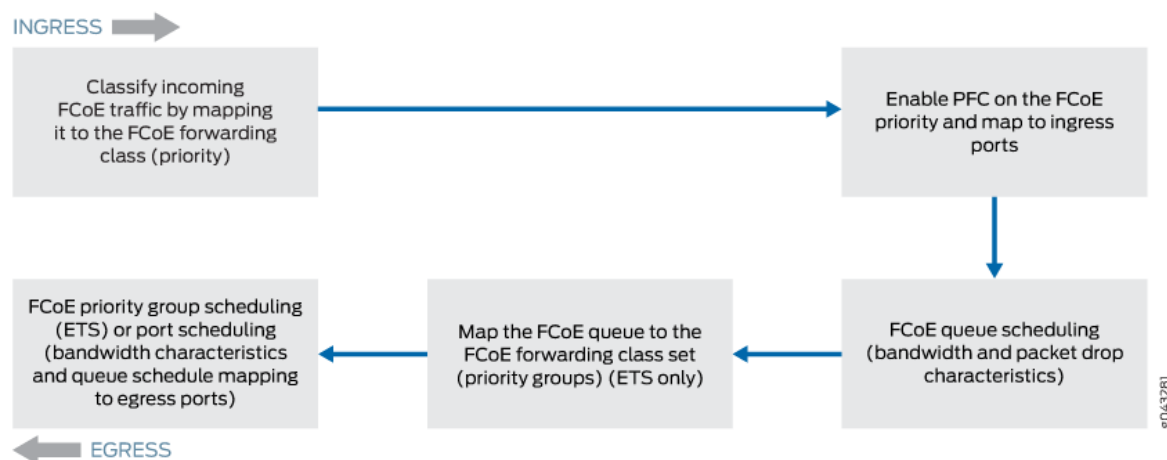
Component	Settings
Hardware	One switch

Table 90: Components of the PFC for FCoE Traffic Configuration Topology (Continued)

Component	Settings
Behavior aggregate classifier (maps the FCoE forwarding class to incoming packets by IEEE 802.1 code point)	Code point 011 to forwarding class fcoe and loss priority low Ingress interfaces: xe-0/0/31, xe-0/0/32, xe-0/0/33, xe-0/0/34
PFC congestion notification profile	fcoe-cnp: Code point 011 Ingress interfaces: xe-0/0/31, xe-0/0/32, xe-0/0/33, xe-0/0/34
FCoE queue scheduler	fcoe-sched: Minimum bandwidth 3g Maximum bandwidth 100% Priority low
Forwarding class-to-scheduler mapping	Scheduler map fcoe-map: Forwarding class fcoe Scheduler fcoe-sched On switches that support direct port scheduling, if you use port scheduling, attach the scheduler map directly to interfaces xe-0/0/31, xe-0/0/32, xe-0/0/33, and xe-0/0/34.
ETS only: Forwarding class set (FCoE priority group)	fcoe-pg: Forwarding class fcoe Egress interfaces: xe-0/0/31, xe-0/0/32, xe-0/0/33, xe-0/0/34
ETS only: Traffic control profile	fcoe-tcp: Scheduler map fcoe-map Minimum bandwidth 3g Maximum bandwidth 100% For ETS hierarchical scheduling, attach the traffic control profile (using the output-traffic-control-profile keyword) to interfaces xe-0/0/31, xe-0/0/32, xe-0/0/33, and xe-0/0/34.

Figure 22 on page 514 shows a block diagram of the configuration components and the configuration flow of the CLI statements used in the example.

Figure 22: PFC for FCoE Traffic Configuration Components Block Diagram



Configuration

IN THIS SECTION

- [CLI Quick Configuration | 514](#)
- [Common Configuration \(Applies to ETS Hierarchical Scheduling and to Port Scheduling\) | 516](#)
- [ETS Hierarchical Scheduling Configuration | 517](#)
- [Port Scheduling Configuration | 518](#)
- [Results | 518](#)

CLI Quick Configuration

To quickly configure PFC for FCoE traffic, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level.

The configuration is separated into the configuration common to ETS and direct port scheduling, and the portions of the configuration that apply only to ETS and only to port scheduling.

Common Configuration that applies to ETS Hierarchical Scheduling and to Port Scheduling:

```
[edit class-of-service]
set classifiers ieee-802.1 fcoe-classifier forwarding-class fcoe loss-priority low code-points 011
set congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
set interfaces xe-0/0/31 unit 0 classifiers ieee-802.1 fcoe-classifier
set interfaces xe-0/0/32 unit 0 classifiers ieee-802.1 fcoe-classifier
set interfaces xe-0/0/33 unit 0 classifiers ieee-802.1 fcoe-classifier
set interfaces xe-0/0/34 unit 0 classifiers ieee-802.1 fcoe-classifier
set interfaces xe-0/0/31 congestion-notification-profile fcoe-cnp
set interfaces xe-0/0/32 congestion-notification-profile fcoe-cnp
set interfaces xe-0/0/33 congestion-notification-profile fcoe-cnp
set interfaces xe-0/0/34 congestion-notification-profile fcoe-cnp
set schedulers fcoe-sched priority low transmit-rate 3g
set schedulers fcoe-sched shaping-rate percent 100
set scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
```

Configuration for ETS hierarchical scheduling—the ETS-specific portion of this example configures forwarding class set (priority group) membership, priority group CoS settings (traffic control profile), and assigns the priority group and its CoS configuration to the interfaces:

```
[edit class-of-service]
set forwarding-class-sets fcoe-pg class fcoe
set traffic-control-profiles fcoe-tcp scheduler-map fcoe-map guaranteed-rate 3g
set traffic-control-profiles fcoe-tcp shaping-rate percent 100
set interfaces xe-0/0/31 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set interfaces xe-0/0/32 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set interfaces xe-0/0/33 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set interfaces xe-0/0/34 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
```

Configuration for port scheduling—the port-scheduling-specific portion of this example assigns the scheduler map (which sets the CoS treatment of the forwarding classes in the scheduler map) to the interfaces:

```
[edit class-of-service]
set interfaces xe-0/0/31 scheduler-map fcoe-map
set interfaces xe-0/0/32 scheduler-map fcoe-map
set interfaces xe-0/0/33 scheduler-map fcoe-map
set interfaces xe-0/0/34 scheduler-map fcoe-map
```

Common Configuration (Applies to ETS Hierarchical Scheduling and to Port Scheduling)

Step-by-Step Procedure

To configure the ingress classifier for FCoE traffic, PFC on the FCoE traffic, apply the PFC and classifier configurations to interfaces, and configure queue scheduling, for both ETS hierarchical scheduling and port scheduling (common configuration):

1. Configure a classifier to set the loss priority and IEEE 802.1 code point assigned to the FCoE forwarding class at the ingress:

```
[edit class-of-service]
user@switch# set classifiers ieee-802.1 fcoe-classifier forwarding-class fcoe loss-priority
low code-points 011
```

2. Configure PFC on the FCoE queue by applying FCoE to the IEEE 802.1 code point 011:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
```

3. Apply the PFC configuration to the ingress interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/31 congestion-notification-profile fcoe-cnp
user@switch# set interfaces xe-0/0/32 congestion-notification-profile fcoe-cnp
user@switch# set interfaces xe-0/0/33 congestion-notification-profile fcoe-cnp
user@switch# set interfaces xe-0/0/34 congestion-notification-profile fcoe-cnp
```

4. Assign the classifier to the ingress interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/31 unit 0 classifiers ieee-802.1 fcoe-classifier
user@switch# set interfaces xe-0/0/32 unit 0 classifiers ieee-802.1 fcoe-classifier
user@switch# set interfaces xe-0/0/33 unit 0 classifiers ieee-802.1 fcoe-classifier
user@switch# set interfaces xe-0/0/34 unit 0 classifiers ieee-802.1 fcoe-classifier
```

5. Configure output scheduling for the FCoE queue:

```
[edit class-of-service]
user@switch# set schedulers fcoe-sched priority low transmit-rate 3g
user@switch# set schedulers fcoe-sched shaping-rate percent 100
```

6. Map the FCoE forwarding class to the FCoE scheduler:

```
[edit class-of-service]
user@switch# set scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
```

ETS Hierarchical Scheduling Configuration

Step-by-Step Procedure

To configure the forwarding class set (priority group) and priority group scheduling (in a traffic control profile), and apply the ETS hierarchical scheduling for FCoE traffic to interfaces:

1. Configure the forwarding class set for the FCoE traffic:

```
[edit class-of-service]
user@switch# set forwarding-class-sets fcoe-pg class fcoe
```

2. Define the traffic control profile for the FCoE forwarding class set:

```
[edit class-of-service]
user@switch# set traffic-control-profiles fcoe-tcp scheduler-map fcoe-map guaranteed-rate 3g
user@switch# set traffic-control-profiles fcoe-tcp shaping-rate percent 100
```

3. Apply the FCoE forwarding class set and traffic control profile to the egress ports:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/31 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
user@switch# set interfaces xe-0/0/32 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
```



```

user@switch# set interfaces xe-0/0/33 forwarding-class-set fcoe-pg output-traffic-control-
profile fcoe-tcp
user@switch# set interfaces xe-0/0/34 forwarding-class-set fcoe-pg output-traffic-control-
profile fcoe-tcp

```

Port Scheduling Configuration

Step-by-Step Procedure

To apply port scheduling for FCoE traffic to interfaces:

1. Apply the scheduler map to the egress ports:

```

[edit class-of-service]
user@switch# set interfaces xe-0/0/31 scheduler-map fcoe-map
user@switch# set interfaces xe-0/0/32 scheduler-map fcoe-map
user@switch# set interfaces xe-0/0/33 scheduler-map fcoe-map
user@switch# set interfaces xe-0/0/34 scheduler-map fcoe-map

```

Results

Display the results of the configuration (the system shows only the explicitly configured parameters; it does not show default parameters such as the fcoe lossless forwarding class). The results are from the ETS hierarchical scheduling configuration to show the more complex configuration. Direct port scheduling results would not show the traffic control profile or forwarding class set portions of the configuration, and would display the name of the scheduler map under each interface (instead of the names of the forwarding class set and output traffic control profile), but is otherwise the same.

```

user@switch> show configuration class-of-service
classifiers {
  ieee-802.1 fcoe-classifier {
    forwarding-class fcoe {
      loss-priority low code-points 011;
    }
  }
}
traffic-control-profiles {
  fcoe-tcp {
    scheduler-map fcoe-map;
    shaping-rate percent 100;
    guaranteed-rate 3000000000;
  }
}

```

```

    }
}
forwarding-class-sets {
    fcoe-pg {
        class fcoe;
    }
}
congestion-notification-profile {
    fcoe-cnp {
        input {
            ieee-802.1 {
                code-point 011 {
                    pfc;
                }
            }
        }
    }
}
}
interfaces {
    xe-0/0/31 {
        congestion-notification-profile fcoe-cnp;
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        unit 0 {
            classifiers {
                ieee-802.1 fcoe-classifier;
            }
        }
    }
    xe-0/0/32 {
        congestion-notification-profile fcoe-cnp;
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        unit 0 {
            classifiers {
                ieee-802.1 fcoe-classifier;
            }
        }
    }
}

```

```

    }
}
xe-0/0/33 {
    congestion-notification-profile fcoe-cnp;
    forwarding-class-set {
        fcoe-pg {
            output-traffic-control-profile fcoe-tcp;
        }
    }
    unit 0 {
        classifiers {
            ieee-802.1 fcoe-classifier;
        }
    }
}
xe-0/0/34 {
    congestion-notification-profile fcoe-cnp;
    forwarding-class-set {
        fcoe-pg {
            output-traffic-control-profile fcoe-tcp;
        }
    }
    unit 0 {
        classifiers {
            ieee-802.1 fcoe-classifier;
        }
    }
}
}
scheduler-maps {
    fcoe-map {
        forwarding-class fcoe scheduler fcoe-sched;
    }
}
schedulers {
    fcoe-sched {
        transmit-rate 3000000000;
        shaping-rate percent 100;
        priority low;
    }
}
}

```

TIP: To quickly configure the interfaces, issue the `load merge` terminal command and then copy the hierarchy and paste it into the switch terminal window.

Verification

IN THIS SECTION

- [Verifying That Priority-Based Flow Control Has Been Enabled | 521](#)
- [Verifying the Ingress Interface PFC Configuration | 522](#)

To verify that the PFC configuration for FCoE traffic components has been created and is operating properly, perform these tasks:

Verifying That Priority-Based Flow Control Has Been Enabled

Purpose

Verify that PFC is enabled on the FCoE queue to enable lossless transport.

Action

List the congestion notification profiles using the operational mode command `show class-of-service congestion-notification`:

```
user@switch> show class-of-service congestion-notification
Type: Input, Name: fcoe-cnp, Index: 51697
Cable Length: 100 m
  Priority    PFC        MRU
  -----
  000        Disabled
  001        Disabled
  010        Disabled
  011        Enabled   2500
  100        Disabled
  101        Disabled
  110        Disabled
  111        Disabled
```

Type: Output	
Priority	Flow-Control-Queues
000	
	0
001	
	1
010	
	2
011	
	3
100	
	4
101	
	5
110	
	6
111	
	7

Meaning

The `show class-of-service congestion-notification` operational command lists all of the congestion notification profiles and which IEEE 802.1p code points have PFC enabled. The command output shows that PFC is enabled on code point 011 for the `fcoe-cnp` congestion notification profile.

The command also shows the default cable length (100 meters), the default maximum receive unit (2500 bytes), and the default mapping of priorities to output queues because this example does not include configuring these options.

Verifying the Ingress Interface PFC Configuration

Purpose

Verify that the classifier `fcoe-classifier` and the congestion notification profile `fcoe-cnp` are configured on ingress interfaces `xe-0/0/31`, `xe-0/0/32`, `xe-0/0/33`, and `xe-0/0/34`.

Action

List the ingress interfaces using the operational mode command `show configuration class-of-service interfaces`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/31
congestion-notification-profile fcoe-cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe-classifier;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/32
congestion-notification-profile fcoe-cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe-classifier;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/33
congestion-notification-profile fcoe-cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe-classifier;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/34
congestion-notification-profile fcoe-cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe-classifier;
    }
}
```

Meaning

The `show configuration class-of-service interfaces` commands list the congestion notification profile that is mapped to the interface (`fcoe-cnp`) and the IEEE 802.1p classifier associated with the interface (`fcoe-classifier`).

RELATED DOCUMENTATION

| *Understanding CoS Flow Control (Ethernet PAUSE and PFC)*

Troubleshooting Dropped FCoE Traffic

IN THIS SECTION

- [Problem | 524](#)
- [Cause | 524](#)
- [Solution | 525](#)

Problem

Description

Fibre Channel over Ethernet (FCoE) traffic for which you want guaranteed delivery is dropped.

Cause

There are several possible causes of dropped FCoE traffic (the list numbers of the possible causes correspond to the list numbers of the solutions in the *Solution* section.):

1. Priority-based flow control (PFC) is not enabled on the FCoE priority (IEEE 802.1p code point) in both the input and output stanzas of the congestion notification profile.
2. The FCoE traffic is not classified correctly at the ingress interface. FCoE traffic should either use the default `fcoe` forwarding class and classifier configuration (maps the `fcoe` forwarding class to IEEE 802.1p code point 011) or be mapped to a lossless forwarding class and to the code point enabled for PFC on the input and output interfaces.

3. The congestion notification profile that enables PFC on the FCoE priority is not attached to the interface.
4. The forwarding class set (priority group) used for guaranteed delivery traffic does not include the forwarding class used for FCoE traffic.

NOTE: This issue can occur only on switches that support enhanced transmission selection (ETS) hierarchical port scheduling. (Direct port scheduling does not use forwarding class sets.)

5. Insufficient bandwidth has been allocated for the FCoE queue or for the forwarding class set to which the FCoE queue belongs.

NOTE: This issue can occur for forwarding class sets only on switches that support ETS hierarchical port scheduling. (Direct port scheduling does not use forwarding class sets.)

6. If you are using Junos OS Release 12.2, the `fcoe` forwarding class has been explicitly configured instead of using the default `fcoe` forwarding class configuration (forwarding-class-to-queue mapping).

NOTE: If you are using Junos OS Release 12.2, use the default forwarding-class-to-queue mapping for the lossless `fcoe` and `no-loss` forwarding classes. If you explicitly configure the lossless forwarding classes, the traffic mapped to those forwarding classes is treated as lossy (best effort) traffic and does *not* receive lossless treatment.

7. If you are using Junos OS Release 12.3 or later and you are not using the default `fcoe` forwarding class configuration, the forwarding class used for FCoE is not configured with the `no-loss` packet drop attribute. In Junos OS 12.3 or later, explicit forwarding classes configurations must include the `no-loss` packet drop attribute to be treated as lossless forwarding classes.

Solution

The list numbers of the possible solutions correspond to the list numbers of the causes in the *Cause* section.

1. Check the congestion notification profile (CNP) to see if PFC is enabled on the FCoE priority (the correct IEEE 802.1p code point) on both input and output interfaces. Use the `show class-of-service congestion-notification` operational command to show the code points that are enabled for PFC in each CNP.

If you are using the default configuration, FCoE traffic is mapped to code point 011 (priority 3). In this case, the input stanza of the CNP should show that PFC is enabled on code point 011, and the output stanza should show that priority 011 is mapped to flow control queue 3.

If you explicitly configured a forwarding class for FCoE traffic, ensure that:

- You specified the `no-loss` packet drop attribute in the forwarding class configuration
- The code point mapped to the FCoE forwarding class in the ingress classifier is the code point enabled for PFC in the CNP input stanza
- The code point and output queue used for FCoE traffic are mapped to each other in the CNP output stanza (if you are not using the default priority and queue, you must explicitly configure each output queue that you want to respond to PFC messages)

For example, if you explicitly configure a forwarding class for FCoE traffic that is mapped to output queue 5 and to code point 101 (priority 5), the output of the `show class-of-service congestion-notification` looks like:

```
Name: fcoe_p5_cnp, Index: 12183
Type: Input
Cable Length: 100 m
  Priority    PFC          MRU
  000        Disabled
  001        Disabled
  010        Disabled
  011        Disabled
  100        Disabled
  101        Enabled    2500
  110        Disabled
  111        Disabled
Type: Output
  Priority    Flow-Control-Queues
  101
           5
```

2. Use the `show class-of-service classifier type ieee-802.1p operational` command to check if the classifier maps the forwarding class used for FCoE traffic to the correct IEEE 802.1p code point.
3. Ensure that the congestion notification profile and classifier are attached to the correct ingress interface. Use the operational command `show configuration class-of-service interfaces interface-name`.

4. Check that the forwarding class set includes the forwarding class used for FCoE traffic. Use the operational command `show configuration class-of-service forwarding-class-sets` to show the configured priority groups and their forwarding classes.
5. Verify the amount of bandwidth allocated to the queue mapped to the FCoE forwarding class and to the forwarding class set to which the FCoE traffic queue belongs. Use the `show configuration class-of-service schedulers scheduler-name` operational command (specify the scheduler for FCoE traffic as the *scheduler-name*) to see the minimum guaranteed bandwidth (transmit-rate) and maximum bandwidth (shaping-rate) for the queue.

Use the `show configuration class-of-service traffic-control-profiles traffic-control-profile` operational command (specify the traffic control profile used for FCoE traffic as the *traffic-control-profile*) to see the minimum guaranteed bandwidth (guaranteed-rate) and maximum bandwidth (shaping-rate) for the forwarding class set.

6. Delete the explicit FCoE forwarding-class-to-queue mapping so that the system uses the default FCoE forwarding-class-to-queue mapping. Include the `delete forwarding-classes class fcoe queue-num 3` statement at the `[edit class-of-service]` hierarchy level to remove the explicit configuration. The system then uses the default configuration for the FCoE forwarding class and preserves the lossless treatment of FCoE traffic.
7. Use the `show class-of-service forwarding-class` operational command to display the configured forwarding classes. The *No-Loss* column shows whether lossless transport is enabled or disabled for each forwarding class. If the forwarding class used for FCoE traffic is not enabled for lossless transport, include the `no-loss` packet drop attribute in the forwarding class configuration (`set class-of-service forwarding-classes class fcoe-forwarding-class-name queue-num queue-number no-loss`).

See [Example: Configuring CoS PFC for FCoE Traffic](#) for step-by-step instructions on how to configure PFC for FCoE traffic, including classifier, interface, congestion notification profile, PFC, and bandwidth scheduling configuration.

RELATED DOCUMENTATION

[show class-of-service congestion-notification](#)

[Configuring CoS PFC \(Congestion Notification Profiles\)](#)

Example: Configuring CoS PFC for FCoE Traffic

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows

IN THIS SECTION

- [Default Lossless Priority Configuration | 529](#)
- [Configuring Lossless Priorities | 531](#)
- [Configuration Rules and Recommendations | 546](#)
- [Lossless Transport Features Introduced in Junos OS Release 12.3 \(Legacy Non-ELS CLI\) | 547](#)
- [Backward Compatibility with Junos OS Releases Earlier Than Release 12.3 \(Legacy Non-ELS CLI\) | 547](#)

The switch supports up to six lossless forwarding classes. (Junos OS Release 12.3 increased support for lossless priorities from two lossless forwarding classes—the default `fcoe` and `no-loss` forwarding classes—to a maximum of six lossless forwarding classes.) Each forwarding class is mapped to an IEEE 802.1p code point (priority).

NOTE: Junos OS Release 13.1 introduced support for up to six lossless forwarding classes on QFabric systems. Throughout this document, features introduced on standalone switches in Junos OS Release 12.3 are introduced on QFabric systems in Junos OS Release 13.1 unless otherwise noted.

Only switches with native Fibre Channel (FC) interfaces, such as the QFX3500, support native FC traffic and configuration as an FCoE-FC gateway. Throughout this document, features that pertain to native FC traffic and to FCoE-FC gateway configuration apply only to switches that support native FC interfaces.



Video: [Why Use PFC in a Data Center Network?](#)

The default configuration is the same as the default configuration in Junos OS Release 12.2 and is backward-compatible. If you need only two (or fewer) lossless forwarding classes, use the default configuration, in which the `fcoe` and `no-loss` forwarding classes are lossless. If you need more than two lossless forwarding classes, you can use the two default lossless forwarding classes and configure additional lossless forwarding classes. If you do not want to use the default lossless forwarding classes, you can change them, or use only the lossless forwarding classes that you explicitly configure.

Default Lossless Priority Configuration

If you do not explicitly configure forwarding classes, the system uses the default forwarding class configuration, which provides two default lossless forwarding classes (*fcoe* and *no-loss*). (If you change the forwarding class configuration, the changes apply to all traffic on that device because forwarding classes are global to a particular device.)

If you do not explicitly configure classifiers, and you do not explicitly configure flow control to pause output queues (configured in the output stanza of the CNP), the default classifier and the default output queue pause configurations are applied to all Ethernet interfaces on the switches (or Node devices). You can override the default classifier and the default output queue pause configuration on a per-interface basis by applying an explicit configuration to an Ethernet interface. The default configuration is used on all Ethernet interfaces that do not have an explicit configuration.

NOTE: If you do not configure flow control on output queues, the default configuration uses a one-to-one mapping of IEEE 802.1p code points (priorities) to output queues by number. For example, priority 0 (code point 000) is mapped to queue 0, priority 1 (code point 001) is mapped to queue 1, and so on. If you do not use the default configuration, you must explicitly configure flow control on each output queue that you want to enable for PFC pause in the output stanza of the CNP.

In the default configuration, only queue 3 and queue 4 are enabled to respond to pause messages from the connected peer. For queue 3 to respond to pause messages, priority 3 (code point 011) must be enabled for PFC in the input stanza of the CNP. For queue 4 to respond to pause messages, priority 4 (code point 100) must be enabled for PFC in the input stanza of the CNP.

The default configuration provides the following lossless behavior:

- Two default lossless forwarding classes (the *no-loss* packet drop attribute is applied to these forwarding classes automatically):
fcoe—Mapped to output queue 3
no-loss—Mapped to output queue 4
- A default classifier that maps the *fcoe* forwarding class to IEEE 802.1p priority 3 (011) and the *no-loss* forwarding class to IEEE 802.1p priority 4 (100)
- Priority-based flow control (PFC) enabled on Ethernet interface output queues 3 and 4 when those queues carry lossless traffic (traffic that is mapped to the *fcoe* and *no-loss* forwarding classes, respectively).

On switches that can be configured as an FCoE-FC gateway, native FC interfaces (NP_Ports), with default flow control enabled on output queue 3 (IEEE 802.1p priority 3) for FCoE/FC traffic.

- DCBX is enabled on all interfaces in autonegotiation mode, and automatically exchanges FCoE application protocol type, length, and values (TLVs) on interfaces that carry FCoE traffic. However, if you explicitly configure DCBX protocol TLV exchange for any application, then you must explicitly configure protocol TLV exchange for every application for which you want DCBX to exchange TLVs, including FCoE.
- On Ethernet ports, PFC buffer calculations use the following default values to determine the headroom buffer size:
Cable length—100 meters (approximately 328 feet)
MRU for priority 3 traffic—2500 bytes
MRU for priority 4 traffic—9216 bytes
Maximum transmission unit (MTU)—1522 (or the configured MTU value for the interface)

NOTE: If you configure flow control on a priority that is not one of the default flow control priorities, the default MRU value is 2500 bytes. For example, if you configure flow control on priority 5 and you do not configure an MRU value, the default MRU value is 2500 bytes.

NOTE: In addition, to support lossless transport, PFC must be enabled explicitly on the lossless IEEE 802.1p priorities (code points) on ingress Ethernet interfaces; no default PFC configuration is applied at ingress interfaces. If you do not enable PFC on lossless priorities, those priorities might experience packet loss during periods of congestion. For example, if you want lossless FCoE traffic and you are using the default fcoe forwarding class, you use a CNP to enable PFC on priority 3 (code point 011), and apply that CNP to all ingress interfaces that carry FCoE traffic.

You can override the default classifier and the default output queue pause configuration on a per-interface basis by applying an explicit configuration to an Ethernet interface.

The default CoS configuration is backward-compatible with the *default* CoS configuration of software releases before Junos OS Release 12.3. If you explicitly configure lossless transport, ensure that the input and output queues corresponding to the lossless forwarding classes are explicitly configured for PFC pause.

[Table 91 on page 531](#) summarizes the default forwarding classes and their mapping to output queues, IEEE 802.1p priorities, and drop attributes.

Table 91: Mapping of Default Forwarding Class to Queue, IEEE 802.1p Priority, and Drop Attribute

Forwarding Class Name	Output Queue	Priority	Drop Attribute
best-effort	0	0	drop
fcoe	3	3	no-loss
no-loss	4	4	no-loss
network-control	7	7	drop

On switches that use the same forwarding classes and output queues for unicast and multdestination (multicast, broadcast, and destination lookup fail) traffic, these forwarding classes carry both unicast and multdestination traffic. Only unicast traffic is treated as lossless traffic. Multdestination traffic is not treated as lossless traffic, even on lossless output queues.

On switches that use different forwarding classes and output queues for unicast and multdestination traffic, there is one default multdestination forwarding class named *mcast*, which is mapped to output queue 8 with a drop attribute of drop. (Incoming multdestination traffic on all IEEE 802.1p priorities is mapped to the mcast forwarding class by default.)

Configuring Lossless Priorities

To configure more than two lossless priorities (forwarding classes), or to change the default mapping of lossless forwarding classes to priorities and paused output queues, you must explicitly configure the switch instead of using the default configuration. Configuring lossless priorities includes:

- Configuring forwarding classes with the no-loss packet drop attribute.
- Using a CNP to configure PFC on ingress interfaces and flow control (PFC) on egress interfaces.
- Configuring a classifier to map IEEE 802.1p priorities (code points) to the correct forwarding classes (the forwarding classes for which you want lossless transport).

NOTE: If you expect a large amount of lossless traffic on your network and configure multiple lossless traffic classes, ensure that you reserve enough scheduling resources (bandwidth) and buffer space to support the lossless flows. (For switches that support shared buffer configuration, [Understanding CoS Buffer Configuration](#) describes how to configure buffers and

provides a recommended buffer configuration for networks with larger amounts of lossless traffic. Buffer optimization is automatic on switches that use virtual output queues.)

In addition, on Ethernet interfaces, DCBX must exchange the appropriate application protocol TLVs for the lossless traffic. On switches that can act as an FCoE-FC gateway, you need to remap the FCoE priority on native FC interfaces if your network uses a priority other than 3 (IEEE code point 011) for FCoE traffic. This section describes:

Configuring Lossless Forwarding Classes (Packet Drop Attribute)

Junos OS Release 12.3 introduced the *no-loss* parameter for forwarding class configuration. (Although it uses the same name, this is not the no-loss default forwarding class. It is a packet drop attribute you can specify to configure any forwarding class as a lossless forwarding class.)

NOTE: On switches that use different forwarding classes for unicast and multdestination traffic, the forwarding class must be a unicast forwarding class. On switches that use the same forwarding classes for unicast and multdestination traffic, only unicast traffic receives lossless treatment.

You can configure up to six forwarding classes (depending on system architecture and the availability of system resources) as lossless forwarding classes by including the *no-loss* drop attribute at the [edit class-of-service forwarding-classes class *forwarding-class-name* queue-num *queue-number*] hierarchy level.

If you use the default *fcoe* or *no-loss* forwarding classes, they include the *no-loss* drop attribute by default. If you explicitly configure the *fcoe* or *no-loss* forwarding classes and you want to retain their lossless behavior, you *must* include the *no-loss* drop attribute in the configuration.

NOTE: All forwarding classes mapped to the same output queue must have the same packet drop attribute. (All forwarding classes mapped to the same output queue must be either lossy or lossless. You cannot map both a lossy and a lossless forwarding class to the same queue.)

To avoid fate sharing (a congested flow affecting an uncongested flow), use a one-to-one mapping of lossless forwarding classes to IEEE 802.1p code points (priorities) and queues. Map each lossless forwarding class to a different queue, and classify incoming traffic into forwarding classes so that each forwarding class transports traffic of only one priority (code point).

The *fcoe* and *no-loss* forwarding classes are special cases, because in the default configuration, they are configured for lossless behavior (providing that you also enable PFC on the priorities mapped to the *fcoe* and *no-loss* forwarding classes in the CNP input stanza).

Table 92 on page 533 summarizes the possible configurations of the fcoe and no-loss forwarding classes in Junos OS Release 12.3 and later, and the result of those configurations in terms of lossless traffic behavior. It is assumed that PFC, DCBX, and classifiers are properly configured.

Table 92: FCoE and No-Loss Forwarding Class Configuration in Junos OS Release 12.3

Explicit (User-Configured) or Default Forwarding Class Configuration	Packet Drop Attribute	Result and Notes
Default	Default	The fcoe and no-loss forwarding classes are lossless. NOTE: Even if you explicitly configure other forwarding classes (lossy or lossless forwarding classes), the fcoe and no-loss forwarding classes remain lossless because they are not explicitly configured.
Explicit	Not specified in the explicit forwarding class configuration	The fcoe and no-loss forwarding classes are lossy because they do not include the no-loss drop attribute.
Explicit	No-loss	The fcoe and no-loss forwarding classes are lossless.
Explicit, configured in Junos OS Release 12.2 or earlier	Not specified (packet drop attribute was not available before Junos OS Release 12.3)	The fcoe and no-loss forwarding classes are lossy in Junos OS Release 12.3 and later because they do not include the no-loss drop attribute. NOTE: To retain lossless behavior, before you upgrade to Junos OS Release 12.3, delete the explicit configuration so that the system uses the default configuration. Alternatively, you can reconfigure the forwarding classes with the no-loss packet drop attribute after upgrading to Junos OS Release 12.3 or later.

For all other forwarding classes except the fcoe and no-loss forwarding classes, you must explicitly configure lossless transport by specifying the no-loss packet drop attribute, because the default configuration for all other forwarding classes is lossy (the no-loss packet drop attribute is not applied).

Congestion Notification Profiles (PFC Configuration)

Use CNPs to configure lossless PFC characteristics on input and output interfaces.

The input stanza of a CNP enables PFC on specified IEEE 802.1p priorities (code points) and fine-tunes headroom buffer settings by configuring the maximum receive unit (MRU) value and cable length on ingress interfaces.

The output stanza of a CNP enables PFC (flow control) on output queues for specified IEEE 802.1p priorities so that the queues can respond to PFC pause messages from the connected peer on the priority of your choice. (By default, output queues 3 and 4 respond to received PFC messages when those queues carry lossless traffic in the fcoe and no-loss forwarding classes, respectively.)

To achieve lossless transport, the priority paused at the ingress interfaces must match the priority paused at the egress interfaces for a given traffic flow. For example, if you configure ingress interfaces to pause traffic tagged with IEEE 802.1p priority 5 (code point 101) and priority 5 traffic is mapped to output queue 5, then you must also configure the corresponding output interfaces to pause priority 5 on queue 5. In addition, the forwarding class mapped to queue 5 must be configured as a lossless forwarding class (using the no-loss drop attribute).



CAUTION: Any change to the PFC configuration on a port temporarily blocks the entire port (not just the priorities affected by the PFC change) so that the port can implement the change, then unblocks the port. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

A change to the PFC configuration means any change to a CNP, including changing the input portion of the CNP (enabling or disabling PFC on a priority, or changing the MRU or cable-length values) or changing the output portion the CNP that enables or disables output flow control on a queue. A PFC configuration change only affects ports that use the changed CNP.

The following actions change the PFC configuration:

- Deleting or disabling a PFC configuration (input or output) in a CNP that is in use on one or more interfaces. For example:
 1. An existing CNP with an input stanza that enables PFC on priorities 3, 5, and 6 is configured on interfaces xe-0/0/20 and xe-0/0/21.
 2. We disable the PFC configuration for priority 6 in the input CNP, and then commit the configuration.
 3. The PFC configuration change causes all traffic on interfaces xe-0/0/20 and xe-0/0/21 to stop until the PFC change has been implemented. When the PFC change has been implemented, traffic resumes.
- Configuring a CNP on an interface. (This changes the PFC state by enabling PFC on one or more priorities.)

- Deleting a CNP from an interface. (This changes the PFC state by disabling PFC on one or more priorities.)

Configuring Input Interface Flow Control (PFC and Headroom Buffer Calculation)

On Ethernet interfaces, the input stanza of the CNP enables PFC on specified priorities so that the ingress interface can send a pause message to the connected peer during periods of congestion. Input CNPs also fine-tune the headroom buffers used for PFC support by allowing you to configure the MRU value and cable length (if you do not want to use the default configuration).

Headroom buffers support lossless transport by storing the traffic that arrives at an interface after the interface sends a PFC flow control message to pause incoming traffic. Until the connected peer receives the flow control message and pauses traffic, the interface continues to receive traffic and must buffer it (and the traffic that is still on the wire after the peer pauses) to prevent packet loss.

The system uses the MRU and the length of the attached physical cable to calculate buffer headroom allocation. The default configuration values are:

- MRU for priority 3 traffic—2500 bytes
- MRU for priority 4 traffic—9216 bytes
- Cable length—100 meters (approximately 328 feet)

NOTE: If you configure flow control on a priority that is not one of the default flow control priorities, the default MRU value is 2500 bytes. For example, if you configure flow control on priority 5 and you do not explicitly configure an MRU value, the default MRU value is 2500 bytes.

You can fine-tune the MRU and the cable length to adjust the size of the headroom buffer on an interface. The switch has a shared global buffer pool and dynamically allocates headroom buffer space to lossless queues as needed.

A lower MRU or a shorter cable length reduces the amount of headroom buffer required on an interface and leaves more headroom buffer space for other interfaces. A higher MRU or a longer cable length increases the amount of headroom buffer space required on an interface and leaves less headroom buffer space for other interfaces.

In many cases, you can better utilize the headroom buffers by reducing the MRU value (for example, an MRU of 2180 is sufficient for most FCoE networks) and by reducing the cable length value if the physical cable is less than 100 meters long.

NOTE: When you configure the headroom buffers by changing the MRU or the cable length, and commit the configuration, the system performs a commit check and rejects the configuration if sufficient headroom buffer space is not available.

However, the system does not perform a commit check but instead returns a syslog error if:

- The buffers are configured on a LAG interface.
- The default classifier is used on the interface (instead of a user-configured classifier).
- The interface has not been created yet.

Configuring Output Interface Flow Control (PFC)

On Ethernet interfaces, you can use the output stanza of the CNP to configure flow control on output queues and enable PFC pause response on specified IEEE 802.1p priorities.

NOTE: On switches that use different output queues for unicast and multidestination traffic, the queue must be a unicast output queue.

By default, output queues 3 and 4 are enabled for PFC pause on priorities 3 (IEEE 802.1p code point 011) and 4 (IEEE 802.1p code point 100). The default PFC pause response supports the default lossless forwarding class configuration, which maps the fcoe forwarding class to queue 3 and priority 3, and maps the no-loss forwarding class to queue 4 and priority 4.

Configuring PFC on output queues enables you to pause any priority on any output queue on any Ethernet interface. Output flow control enables you to use more than two output queues to support lossless traffic flows (you can configure up to six lossless forwarding classes and map them to different output queues that are enabled for PFC pause). Output queue flow control also enables you to support multiple lossless forwarding classes (each mapped to a different priority and output queue) for one class of traffic.

NOTE: Output flow control only works when PFC is enabled in the CNP input stanza on the corresponding priorities on the interface. For example, if you enable output flow control on priority 5 (IEEE 802.1p code point 101), then you must also enable PFC in the CNP on the input stanza on priority 5.

For example, if the converged Ethernet network uses two different priorities for FCoE traffic (for example, priority 3 and priority 5), then you can classify those priorities into different lossless forwarding classes that are mapped to different output queues:

1. Configure two lossless forwarding classes for FCoE traffic, with each forwarding class mapped to a different output queue. For example, you could use the default fcoe forwarding class, which is mapped to queue 3, and you could configure a second lossless forwarding class called fcoe1 and map it to queue 5. The fcoe forwarding class is for priority 3 FCoE traffic (code point 011), and the fcoe1 forwarding class is for priority 5 (code point 101) FCoE traffic.
2. Configure a classifier that maps each forwarding class to the desired IEEE 802.1p code point (priority). If FCoE traffic on both priorities uses one interface, the classifier must classify both forwarding classes to the correct priorities. If FCoE traffic of different priorities uses different interfaces, the classifier configuration on each interface must map the correct priority to the corresponding lossless forwarding class.
3. Apply the classifier to the interfaces that carry FCoE traffic. The classifier determines the mapping of forwarding classes to priorities on each interface.

To configure lossless transport for these forwarding classes, you also need to:

- Enable PFC on the two priorities (3 and 5 in this example) at the ingress interfaces in the CNP input stanza.
- Configure PFC on the output queues and priorities for the forwarding classes in the CNP output stanza so that the interface can respond to pause messages received from the connected peer.

NOTE: When you configure the CNP on an interface, all ingress and egress traffic is blocked until the configuration is implemented, then the interface is unblocked and traffic resumes. During the time the interface is blocked, all queues on the interface experience packet loss.

- Configure DCBX to exchange application protocol TLVs on both FCoE priorities.

NOTE: If you do not configure flow control to pause output queues, the default configuration uses a one-to-one mapping of IEEE 802.1p code points (priorities) to output queues by number. For example, priority 0 (code point 000) is mapped to queue 0, priority 1 (code point 001) is mapped to queue 1, and so on. By default, only queues 3 and 4 are enabled to respond to pause messages from the connected peer, and you must explicitly enable PFC on the corresponding priorities in the CNP input stanza to achieve lossless behavior.

If you do not use the default configuration, you must explicitly configure flow control on each output queue that you want to enable for PFC pause. For example, if you explicitly configure flow control on output queue 5, the default configuration is no longer valid, and only output queue 5 is enabled for PFC pause. Output queues 3 and 4 are no longer enabled for PFC pause, so traffic using those queues no longer responds to PFC pause messages even if the

corresponding forwarding class is configured with the no-loss drop attribute. To retain the pause configuration on output queues 3 and 4 and configure flow control on queue 5, you need to explicitly configure flow control on queues 3, 4, and 5.

On switches that use different output queues for unicast and multidestination traffic, you cannot configure flow control to pause a multidestination output queue. You can configure flow control to pause only unicast output queues. On switches that use the same output queues for unicast and multidestination traffic, only unicast traffic receives lossless treatment.

Output Interface Flow Control Profiles

Configuring the CNP output stanza creates an output flow control profile that tells egress ports the queues on which the Ethernet interface should respond to PFC pause messages. Although you can create an unlimited number of CNPs that contain input stanzas only, the number of CNPs that you can configure with output stanzas is limited:

- For standalone switches that are not part of a QFabric system, you can configure up to two output interface flow control profiles. (You can configure up to two CNPs with output stanzas.)
- For QFabric systems, you can configure one output interface flow control profile per Node device. (You can configure one CNP with an output stanza per Node device.)

There are a total of four output flow control profiles.

The system has a default output flow control profile that is applied to all Ethernet interfaces when the CNP attached to the interface has only an input stanza and does not include an output stanza. The default profile responds to PFC pause messages received on queue 3 (for priority 3, for the default fcoe forwarding class) and on queue 4 (for priority 4, for the default no-loss forwarding class), and is effective only if PFC is configured on those priorities in the CNP input stanza.

Additionally, the system has two internal output flow control profiles that it applies automatically to fabric (FTE) ports and to native FC interfaces (NP_Ports). When the switch is not part of a QFabric system, the profile normally used for FTE ports is available for user configuration and provides a second user-configurable profile. (That is why standalone switches have two user-configurable output flow control profiles, but Node devices on a QFabric system have only one user-configurable output flow control profile.)

Because one output CNP can configure PFC pause response on multiple output queues (priorities), one user-configurable output CNP is usually flexible enough to specify the desired PFC response on all programmed interfaces.

NOTE: Each port can use one output flow control profile. You cannot apply more than one profile to one port.

Output flow control profiles can be expressed in table format. For example, [Table 93 on page 539](#) shows the default output flow control profile that pauses priorities 3 and 4 on queues 3 and 4 (remember that PFC must also be enabled on code points 3 and 4 in the CNP input stanza in order for PFC to work):

Table 93: Default Output Flow Control Profile

IEEE 802.1p Priority Specified in Received PFC Frame	Paused Output Queue
0 (000)	—
1 (001)	—
2 (010)	—
3 (011)	3
4 (100)	4
5 (101)	—
6 (110)	—
7 (111)	—

[Table 94 on page 540](#) is an example of a user-configured output flow control profile. Using the example from the preceding section, the CNP output stanza configures flow control on output queue 5, and also explicitly configures output flow control on queues 3 and 4 for the fcoe and no-loss forwarding classes. (If you explicitly configure an output CNP, you must explicitly configure every output queue that you want to respond to PFC messages, because the user-configured profile overrides the default profile. If this example did not include queues 3 and 4, those queues would no longer respond to received PFC messages.)

Table 94: User-Configured Output Flow Control Profile

IEEE 802.1p Priority Specified in Received PFC Frame	Paused Output Queue
0 (000)	—
1 (001)	—
2 (010)	—
3 (011)	3
4 (100)	4
5 (101)	5
6 (110)	—
7 (111)	—

Remember that you must also enable PFC on code points 3, 4, and 5 in the CNP input stanza for this configuration to work. When you configure the CNP on an interface, all ingress and egress traffic is blocked until the configuration is implemented, then the interface is unblocked and traffic resumes. During the time the interface is blocked, all queues on the interface experience packet loss.

Configuring PFC Across Layer 3 Interfaces on QFX5210, QFX5200, QFX5100, EX4600, and QFX10000 Switches

Enabling PFC on traffic flows is based on the IEEE 802.1p code point (priority) in the priority code point (PCP) field of the Ethernet frame header (sometimes known as the CoS bits). To enable PFC on traffic that crosses Layer 3 interfaces, the traffic must be classified by its IEEE 802.1p code point, not by its DSCP (or DSCP IPv6) code point.

See [Understanding PFC Functionality Across Layer 3 Interfaces](#) for a conceptual overview of how to enable PFC on traffic across Layer 3 interfaces. See [Example: Configuring PFC Across Layer 3 Interfaces](#) for an example of how to configure PFC on traffic that traverses Layer 3 interfaces.

Configuring DCBX (Application Protocol TLV Exchange)

For applications that require lossless transport, DCBX exchanges application protocol TLVs with the connected peer interface. By default, DCBX advertises FCoE application protocol TLVs on all interfaces that are enabled for DCBX, and by default, DCBX is enabled on all interfaces. DCBX advertises no other applications by default.

For each application (for example, iSCSI) that you want to configure for lossless transport, you must enable the interfaces which carry that application traffic to exchange DCBX protocol TLVs with the connected peer. The TLV exchange allows the peer interfaces to negotiate a compatible configuration to support the application.

If you configure DCBX to advertise any application, the default DCBX advertisement is overridden, and DCBX advertises only the configured applications. If you want an interface to advertise only the FCoE application, you do not have to configure DCBX application protocol TLV exchange; instead, you can use the default configuration.

If you want DCBX to advertise other applications, you must explicitly configure an application map and apply it to the interfaces on which you want to exchange protocol TLVs for those applications. If you want to exchange FCoE application protocol TLVs in addition to other application protocol TLVs, you must also explicitly configure the FCoE application in the application map. [Understanding DCBX Application Protocol TLV Exchange](#) describes how application mapping works.

NOTE: Lossless transport also requires that you enable PFC on the correct priority (IEEE 802.1p code point) on the ingress interfaces using an input CNP. If the priority you pause at the ingress interfaces is not mapped to queue 3 or queue 4 (the two output queues that are enabled for PFC pause flow control by default), then you must also enable the output queues that correspond to paused input priorities to pause using the output stanza of the CNP.

Fate Sharing Among Traffic Classes

You can configure different lossless (or lossy) traffic flows to share fate—that is, to receive the same CoS treatment.

Fate sharing is not desirable for I/O convergence. Instead of independent control of the fate of each type of flow, different types of flows receive the same treatment. Fate sharing is particularly undesirable for lossless flows. If one lossless flow experiences congestion and must be paused, that affects flows that share fate with the congested flow even if the other flows are not experiencing congestion, and also can cause ingress port congestion. If your network requires that all 802.1p priorities be lossless, you can achieve that by allowing some fate sharing among the eight priorities by spreading them across up to six lossless forwarding classes.

If the number of lossless priorities is less than or equal to the number of configured lossless forwarding classes, then you can avoid fate sharing by configuring a one-to-one mapping of forwarding classes to IEEE 802.1p code points (priorities) and output queues. (Each forwarding class should be mapped to a different output queue and classified to a different priority.)

If you want to configure different traffic flows to share fate, two fate-sharing configurations are supported: mapping one forwarding class to more than one IEEE 802.1p code point (priority), and mapping two forwarding classes to the same output queue:

1. If you map one lossless forwarding class to more than one priority, the traffic tagged with each of the priorities uses the same CoS properties associated (the CoS properties associated with the forwarding class). For example, configuring a forwarding class called fc1, mapping it to queue 1, and mapping it to code points 101 and 110 using a classifier named classify1 results in the traffic tagged with priorities 101 and 110 sharing fate:

```
user@switch# set class-of-service forwarding-classes class fc1 queue-num 1 no-loss
user@switch# set class-of-service classifiers ieee-802.1 classify1 forwarding class fc1 loss-
priority low code-points 101
user@switch# set class-of-service classifiers ieee-802.1 classify1 forwarding class fc1 loss-
priority low code-points 110
```

In this case, if the traffic mapped to either priority experiences congestion, both priorities are paused because they are mapped to the same forwarding class and are therefore treated similarly.

2. If you map multiple lossless forwarding classes to the same output queue, the traffic mapped to the forwarding classes uses the same output queue. This increases the amount of traffic on the queue, and can create congestion that affects all of the traffic flows that are mapped to the queue. For example, configuring two forwarding classes called fc1 and fc2, mapping both forwarding classes to queue 1, and mapping the forwarding classes to code points 101 and 110 (respectively) using a classifier named classify1, results in the traffic tagged with priorities 101 and 110 sharing fate on the same output queue:

```
user@switch# set class-of-service forwarding-classes class fc1 queue-num 1 no-loss
user@switch# set class-of-service forwarding-classes class fc2 queue-num 1 no-loss
user@switch# set class-of-service classifiers ieee-802.1 classify1 forwarding class fc1 loss-
priority low code-points 101
user@switch# set class-of-service classifiers ieee-802.1 classify1 forwarding class fc2 loss-
priority low code-points 110
```

In this case, even though the two forwarding classes use different IEEE 802.1p priorities, if one forwarding class experiences congestion, it affects the other forwarding class. The reason is that if

the output queue is paused because of congestion on either forwarding class, all traffic that uses that queue is paused. Since both forwarding classes are mapped to the queue, the traffic mapped to both forwarding classes is paused.

NOTE: If you map more than one forwarding class to a queue, all of the forwarding classes mapped to the same queue must have the same packet drop attribute (all of the forwarding classes must be lossy, or all of the forwarding classes mapped to a queue must be lossless).

Transit Switch Configuration Versus FCoE-FC Gateway Configuration

On a transit switch (all Ethernet ports, no native FC ports) that forwards FCoE traffic (or other traffic that requires lossless transport across the Ethernet network), the configuration of classifiers, lossless forwarding classes, DCBX, and PFC on ingress and egress interfaces to support lossless transport is as described in this document.

When a switch acts as an FCoE-FC gateway (if native FC interfaces are supported on your switch), the system uses native FC interfaces (NP_Ports) to connect to the FC switch (or FCoE forwarder) at the FC network edge. You cannot apply CNPs or DCBX to native FC interfaces, only to Ethernet interfaces.

On an FCoE-FC gateway, the Ethernet interface configuration of classifiers, DCBX, and PFC is the same as the Ethernet interface configuration on a transit switch. The configuration of lossless forwarding classes is also the same.

However, supporting lossless transport on native FC interfaces requires that you rewrite the IEEE 802.1p priority value *if* your network uses any priority other than 3 (IEEE code point 011) for FCoE traffic. If your network uses priority 3 for FCoE traffic, you can and should use the default configuration on native FC interfaces.

By default, native FC interfaces tag packets with priority 3 when they encapsulate the incoming FC packets in Ethernet. If your FCoE network uses a different priority than 3 for FCoE traffic, you need to rewrite the priority value to the value that your network uses on the FC interface, classify the FCoE traffic to the correct priority on the Ethernet interfaces, and enable PFC on the correct priority on the Ethernet interfaces, as described in "[Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway](#)" on page 620.

Configuration Results and Commit Checks

Different configurations of forwarding classes and their drop attributes, classifiers, CNPs (PFC flow control), and Ethernet PAUSE (IEEE 802.3X flow control) result in different system behaviors.

[Table 95 on page 544](#) describes the results of the possible lossless transport configurations in each case. The assumption in the *Result* column is that the system's buffer headroom calculation resulted in a successful configuration.

However, if the system calculates that there is insufficient buffer space to support the configuration, a commit check prevents you from committing the configuration on an individual Ethernet interface. For LAG interfaces, the system does not issue a commit check error but instead issues a syslog message.

NOTE: After you configure lossless transport for a LAG interface, be sure to check the syslog messages to confirm that the commit was successful.

Table 95: Results of Lossless Priority Configuration

Classifier Configuration	Congestion Notification Profile Configuration	Ethernet PAUSE (IEEE 802.3X) Configuration	Result
None (default classifier)	None	None	System default configuration. No flows are lossless. To achieve lossless behavior for the default fcoe and no-loss forwarding classes, you must configure an input CNP to enable PFC on their IEEE 802.1p code points (011 and 100 respectively).
Classifier with no lossless forwarding classes	None	None	No lossless traffic flows are configured; all traffic is best effort.
Classifier with at least one lossless forwarding class	None	None	Because no CNP is attached to interfaces, PFC is not enabled on the code point of the lossless traffic and no headroom buffer is allocated to the lossless queue, so packets can drop during periods of congestion. This configuration does not achieve lossless behavior.
None (default classifier)	PFC enabled on the fcoe and no-loss forwarding class code points (priorities)	None	The default classifier classifies traffic into two lossless forwarding classes, fcoe and no-loss. The CNP enables PFC on the priorities mapped to both lossless forwarding classes, resulting in lossless behavior for traffic mapped to the fcoe and no-loss forwarding classes.

Table 95: Results of Lossless Priority Configuration (Continued)

Classifier Configuration	Congestion Notification Profile Configuration	Ethernet PAUSE (IEEE 802.3X) Configuration	Result
None (default classifier)	None	Flow control enabled	The system calculates buffer headroom for the physical link based on the interface MTU and the default cable length. The system does not calculate buffer headroom for individual output queues. Because Ethernet PAUSE is enabled on the link instead of PFC being enabled on the lossless priorities, the entire link is paused during periods of congestion. This configuration results in lossless behavior for all of the forwarding classes on the link, but because all traffic is paused, this can cause greater overall network congestion.
Classifier with at least one lossless forwarding class	PFC enabled on the lossless forwarding class code points (priorities)	None	Headroom buffer allocated only to priorities that are mapped to the lossless forwarding classes and on which PFC is enabled. This configuration achieves lossless behavior for the lossless forwarding classes.
Classifier with no lossless forwarding classes	None	Flow control enabled	The system calculates buffer headroom for the physical link based on the interface MTU and the default cable length, and it pauses all traffic on the link during periods of congestion.
Classifier with at least one lossless forwarding class	None	Flow control enabled	The system calculates buffer headroom for the physical link based on the interface MTU and the default cable length, and it pauses all traffic on the link during periods of congestion.

Table 95: Results of Lossless Priority Configuration (*Continued*)

Classifier Configuration	Congestion Notification Profile Configuration	Ethernet PAUSE (IEEE 802.3X) Configuration	Result
Classifier with at least one lossless forwarding class	PFC enabled on the lossless forwarding class code points (priorities)	Flow control enabled on a <i>different</i> interface than the interface with the CNP	The system checks the available buffer space for both the PFC-enabled priorities and for the other link. If sufficient buffer space is available, the lossless forwarding classes configured with PFC on one interface and also all of the traffic on the link with Ethernet PAUSE enabled achieve lossless behavior.

NOTE: If you attempt to configure both PFC and Ethernet PAUSE on a link, the system returns a commit error. PFC and Ethernet PAUSE are mutually exclusive configurations on an interface.

Configuration Rules and Recommendations

Keep in mind the following configuration rules and recommendations when you configure lossless traffic flows:

- You can configure a maximum of six lossless forwarding classes (forwarding classes with the no-loss packet drop attribute).
- All forwarding classes that you map to the same queue must have the same packet drop attribute (all of the forwarding classes must be lossy, or all of the forwarding classes must be lossless).
- Do not configure weighted random early detection (WRED) on lossless forwarding classes. (Do not associate a drop profile with a forwarding class that has the no-loss packet drop attribute.)
- On switches that use different forwarding classes and output queues for unicast and multdestination traffic, you cannot configure flow control to pause a multdestination output queue. You can configure PFC flow control only to pause unicast output queues.
- On switches that use different forwarding classes and output queues for unicast and multdestination traffic, forwarding classes mapped to multdestination queues (queues 8 through 11) cannot have the no-loss packet drop attribute. (Multdestination forwarding classes cannot be configured as lossless forwarding classes.)

Lossless Transport Features Introduced in Junos OS Release 12.3 (Legacy Non-ELS CLI)

Support for lossless transport introduced in Junos OS Release 12.3 includes:

- Configuring up to six lossless forwarding classes.
- Configuring PFC pause on output queues to program the output queues that can respond to PFC pause messages received from the connected peer. The priorities you pause on output queues must match the priorities on which you enable PFC on the corresponding ingress interfaces. For example, if you program output queues to pause priorities 3 (011) and 5 (101), then you must also enable pause on priorities 3 and 5 on the corresponding ingress interfaces. Configuring flow control on the output queues and enabling PFC on the corresponding input queues allows you to pause up to six priorities (forwarding classes).
- Controlling the headroom buffer on Ethernet interfaces by configuring the maximum receive unit (MRU) size for the traffic mapped to an IEEE 802.1p priority (configured per priority) and the length of the attached cable (configured per interface). The MRU size can range up to full jumbo packet size (9216 bytes).
- Remapping (rewriting) IEEE 802.1p priorities on native Fibre Channel (FC) interfaces when the system is acting as an FCoE-FC gateway. If the Ethernet (FCoE) network uses a different IEEE 802.1p priority than priority 3 (011) for FCoE traffic, then you can use priority remapping to classify FCoE traffic into a lossless forwarding class mapped to that different priority (see "[Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway](#)" on page 620).

Lossless transport still requires configuring previously existing features, including enabling PFC on the lossless priorities on ingress interfaces, and configuring classifiers to classify incoming traffic into lossless forwarding classes based on the IEEE 802.1p priority tag of the packet.

NOTE: If you expect a large amount of lossless traffic on your network and configure multiple lossless traffic classes, ensure that you reserve enough scheduling resources (bandwidth) and lossless headroom buffer space to support the lossless flows. ([Understanding CoS Buffer Configuration](#) describes how to configure buffers and provides a recommended buffer configuration for networks with larger amounts of lossless traffic.)

Backward Compatibility with Junos OS Releases Earlier Than Release 12.3 (Legacy Non-ELS CLI)

The addition of the no-loss packet drop attribute to forwarding class configuration means that when you upgrade from an earlier release to Junos OS Release 12.3, the new software might not preserve the lossless forwarding class configuration of the fcoe and no-loss forwarding classes.

If you used the default forwarding class configuration for the fcoe and no-loss forwarding classes, the CoS configuration is backward-compatible. You do not have to do anything to preserve the lossless behavior of traffic that uses those forwarding classes when you upgrade to Junos OS Release 12.3. (This is because the default configuration of these two forwarding classes includes the no-loss packet drop attribute.)

However, if you explicitly configured the fcoe or the no-loss forwarding class by including the `set forwarding-classes class forwarding-class-name queue-num queue-number` statement at the [edit class-of-service] hierarchy level, then those forwarding classes are no longer lossless, they are lossy. (They are lossy because explicit configuration in releases earlier than Junos OS Release 12.3 did not use the no-loss packet drop attribute.) In Junos OS Release 12.3 and later, you must include the no-loss packet drop attribute in explicit forwarding class configurations to configure a lossless forwarding class.

For example, before Junos OS Release 12.3, the following explicit configuration resulted in a lossless forwarding class:

```
user@switch# set class-of-service forwarding-classes class fcoe queue-num 3
```

However, in Junos OS Release 12.3, this configuration is lossy because it does not include the no-loss packet drop attribute. To preserve lossless behavior, after upgrading to Junos OS Release 12.3, you need to add the no-loss drop attribute:

```
user@switch# set class-of-service forwarding-classes class fcoe queue-num 3 no-loss
```

Alternatively, you can delete the explicit configuration before you upgrade to Junos OS Release 12.3 so that the system uses the default forwarding class, which is lossless:

```
user@switch# delete class-of-service forwarding-classes class fcoe queue-num 3
```

NOTE: The explicit configuration of other forwarding classes does not affect the lossless (or lossy) state of the fcoe and no-loss forwarding classes, because only the fcoe and no-loss forwarding classes were lossless forwarding classes before Junos OS Release 12.3. For example, if you explicitly configured the best-effort forwarding class but you used the default fcoe and no-loss forwarding classes in Junos OS Release 12.2, then when you upgrade to Junos OS Release 12.3, the fcoe and no-loss forwarding classes are still lossless (and the best-effort forwarding classes retains its explicit configuration).

NOTE: To achieve lossless behavior for the traffic belonging to any forwarding class, you must also use a CNP to enable PFC on the IEEE 802.1p priority mapped to the forwarding class and apply the CNP to the relevant interfaces, and ensure that DCBX exchanges the protocol TLVs for the application with the connected peer.

RELATED DOCUMENTATION

Understanding DCBX Application Protocol TLV Exchange

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Understanding PFC Functionality Across Layer 3 Interfaces

Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch)

Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface

Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces

Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI)

Example: Configuring PFC Across Layer 3 Interfaces

Configuring CoS PFC (Congestion Notification Profiles)

Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch)

IN THIS SECTION

- [Requirements | 550](#)
- [Overview | 550](#)
- [Configuration | 553](#)
- [Verification | 556](#)

The default system configuration supports FCoE traffic on priority 3 (IEEE 802.1p code point 011). If the FCoE traffic on your converged Ethernet network uses priority 3, the only user configuration required for lossless transport is to enable PFC on code point 011 on the FCoE ingress interfaces.

However, if your network uses a different priority than 3 for FCoE traffic, you need to configure lossless FCoE transport on that priority. This example shows you how to configure lossless FCoE transport on a converged Ethernet network that uses priority 5 (IEEE 802.1p code point 101) for FCoE traffic instead of using priority 3.

Requirements

This example uses the following hardware and software components:

- One switch used as an FCoE transit switch
- Junos OS Release 12.3 or later for the QFX Series

Overview

IN THIS SECTION

- [Topology | 551](#)

Although FCoE traffic typically uses IEEE 802.1p priority 3 on converged Ethernet networks, some networks use a different priority for FCoE traffic. Regardless of the priority used, FCoE traffic must receive lossless treatment. Supporting lossless behavior for FCoE traffic when your network does not use priority 3 requires configuring:

- A lossless forwarding class for FCoE traffic.
- A behavior aggregate (BA) classifier to map the FCoE forwarding class to the appropriate IEEE 802.1p priority.
- A congestion notification profile (CNP) to enable PFC on the FCoE code point at the interface ingress and to configure flow control on the interface egress. Flow control on the interface egress enables the interface to respond to PFC messages received from the connected peer and pause the correct IEEE 802.1p priority on the correct output queue.

NOTE: Configuring or changing PFC on an interface blocks the entire port until the PFC change is completed. After a PFC change is completed, the port is unblocked and traffic

resumes. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

- A DCBX application and an application map to support DCBX application TLV exchange for the lossless FCoE traffic on the configured FCoE priority. By default, DCBX is enabled on all Ethernet interfaces, but only on priority 3 (IEEE 802.1p code point 011). To support DCBX application TLV exchange when you are not using the default configuration, you must configure all of the applications and map them to interfaces and priorities.

The priorities specified in the BA classifiers, CNP, and DCBX application map must match, or the configuration does not work. You must specify the same lossless FCoE forwarding class in each configuration and use the same IEEE 802.1p code point (priority) so that the FCoE traffic is properly classified into flows and so that those flows receive lossless treatment.

Topology

This example shows how to configure one lossless FCoE traffic class, map it to a priority other than priority 3, and configure flow control to ensure lossless behavior on the interfaces. This example uses two Ethernet interfaces, xe-0/0/25 and xe-0/0/26. The interfaces connect to a converged Ethernet network that uses IEEE 802.1p priority 5 (code point 101) for FCoE traffic.

The configuration on the two interfaces is the same. Both interfaces use the same explicitly configured lossless FCoE forwarding class and the same ingress classifier. Both interfaces enable PFC on priority 5 and enable flow control on the same output queue (which is mapped to the lossless FCoE forwarding class).

[Table 96 on page 551](#) shows the configuration components for this example.

Table 96: Components of the Configuration Topology for FCoE Traffic That Does Not Use Priority 3

Component	Settings
Hardware	One switch

Table 96: Components of the Configuration Topology for FCoE Traffic That Does Not Use Priority 3
(Continued)

Component	Settings
Forwarding class	<p>Name—fcoe1</p> <p>Queue mapping—queue 5</p> <p>Packet drop attribute—no-loss</p> <p>NOTE: A lossless forwarding class can be mapped to any output queue. However, because the fcoe1 forwarding class uses priority 5 in this example, matching that traffic to a forwarding class that uses queue 5 creates a configuration that is logical and easy to map because the priority and the queue are identified by the same number.</p>
BA classifier	<p>Name—fcoe_p5</p> <p>FCoE priority mapping—Forwarding class fcoe1 mapped to code point 101 (IEEE 802.1p priority 5) and a packet loss priority of low.</p>
PFC configuration (CNPs)	<p>CNP name—fcoe_p5_cnp</p> <p>Input CNP code point—101</p> <p>MRU—2240 bytes</p> <p>Cable length—100 meters</p> <p>Output CNP code point—101</p> <p>Output CNP flow control queue—5</p> <p>NOTE: When you apply a CNP with an explicit output queue flow control configuration to an interface, the explicit CNP overwrites the default output CNP. The output queues that are enabled for pause in the default configuration (queues 3 and 4) are not enabled for pause unless they are included in the explicitly configured output CNP.</p>

Table 96: Components of the Configuration Topology for FCoE Traffic That Does Not Use Priority 3
(Continued)

Component	Settings
DCBX application mapping	Application name—fcoe_p5_app Application EtherType—0x8906 Application map name—fcoe_p5_app_map Application map code points—101 NOTE: LLDP and DCBX must be enabled on the interface. By default, LLDP and DCBX are enabled on all Ethernet interfaces.

NOTE: This example does not include scheduling (bandwidth allocation) configuration or the FIP snooping configuration. This example focuses only on the lossless FCoE priority configuration. QFX10000 switches do not support FIP snooping. For this reason, QFX10000 switches cannot be used as FCoE access transit switches. QFX10000 switches can be used as intermediate or aggregation transit switches in the FCoE path, between an FCoE access transit switch that performs FIP snooping and an FCF.

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 553](#)
- [Configuring A Lossless FCoE Forwarding Class On IEEE 802.1p Priority 5 | 554](#)

CLI Quick Configuration

To quickly configure a lossless FCoE forwarding class that uses a different priority than IEEE 802.1p priority 3 for FCoE traffic on an FCoE transit switch, copy the following commands, paste them in a text

file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level.

```
set class-of-service forwarding-classes class fcoe1 queue-num 5 no-loss
set class-of-service classifiers ieee-802.1 fcoe_p5 forwarding-class fcoe1 loss-priority low
code-points 101
set class-of-service interfaces xe-0/0/25 unit 0 classifiers ieee-802.1 fcoe_p5
set class-of-service interfaces xe-0/0/26 unit 0 classifiers ieee-802.1 fcoe_p5
set class-of-service congestion-notification-profile fcoe_p5_cnp input ieee-802.1 code-point 101
pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p5_cnp input cable-length 100
set class-of-service congestion-notification-profile fcoe_p5_cnp output ieee-802.1 code-point
101 pfc flow-control-queue 5
set class-of-service interfaces xe-0/0/25 congestion-notification-profile fcoe_p5_cnp
set class-of-service interfaces xe-0/0/26 congestion-notification-profile fcoe_p5_cnp
set applications application fcoe_p5_app ether-type 0x8906
set policy-options application-maps fcoe_p5_app_map application fcoe_p5_app code-points 101
set protocols dcbx interface xe-0/0/25 application-map fcoe_p5_app_map
set protocols dcbx interface xe-0/0/26 application-map fcoe_p5_app_map
```

Configuring A Lossless FCoE Forwarding Class On IEEE 802.1p Priority 5

Step-by-Step Procedure

To configure a lossless forwarding class for FCoE traffic on IEEE 802.1p priority 5 (code point 101), classify FCoE traffic into the lossless forwarding class, configure a congestion notification profile to enable PFC on the FCoE priority and output queue, and configure DCBX application protocol TLV exchange for traffic on the FCoE priority:

1. Configure the lossless forwarding class (named `fcoe1` and mapped to output queue 5) for FCoE traffic on IEEE 802.1p priority 5:

```
[edit class-of-service]
user@switch# set forwarding-classes class fcoe1 queue-num 5 no-loss
```

2. Configure the ingress classifier (fcoe_p5). The classifier maps the FCoE priority (code point 101) to the lossless FCoE forwarding class fcoe1:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_p5 forwarding-class fcoe1 loss-priority low code-points 101
```

3. Apply the classifier to interfaces xe-0/0/25 and xe-0/0/26:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/25 unit 0 classifiers ieee-802.1 fcoe_p5
user@switch# set interfaces xe-0/0/26 unit 0 classifiers ieee-802.1 fcoe_p5
```

4. Configure the CNP. The input stanza enables PFC on the FCoE priority (IEEE 802.1p code point 101), sets the MRU value (2240 bytes), and sets the cable length value (100 meters). The output stanza configures flow control on output queue 5 on the FCoE priority:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe_p5_cnp input ieee-802.1 code-point 101
pfc mru 2240
user@switch# set congestion-notification-profile fcoe_p5_cnp input cable-length 100
user@switch# set congestion-notification-profile fcoe_p5_cnp output ieee-802.1 code-point 101
pfc flow-control-queue 5
```

5. Apply the CNP to the interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/25 congestion-notification-profile fcoe_p5_cnp
user@switch# set interfaces xe-0/0/26 congestion-notification-profile fcoe_p5_cnp
```

6. Configure the DCBX application for FCoE to map to the Ethernet interfaces, so that DCBX can exchange application protocol TLVs on the IEEE 802.1p priority 5 instead of on the default priority 3:

```
[edit]
user@switch# set applications application fcoe_p5_app ether-type 0x8906
```

7. Configure a DCBX application map to map the FCoE application to the correct IEEE 802.1p FCoE priority:

```
[edit]
user@switch# set policy-options application-maps fcoe_p5_app_map application fcoe_p5_app code-
points 101
```

8. Apply the application map to the Ethernet interfaces so that DCBX exchanges FCoE application TLVs on the correct code point:

```
[edit]
user@switch# set protocols dcbx interface xe-0/0/25 application-map fcoe_p5_app_map
user@switch# set protocols dcbx interface xe-0/0/26 application-map fcoe_p5_app_map
```

Verification

IN THIS SECTION

- [Verifying the Forwarding Class Configuration | 556](#)
- [Verifying the Behavior Aggregate Classifier Configuration | 557](#)
- [Verifying the PFC Flow Control Configuration \(CNP\) | 558](#)
- [Verifying the Interface Configuration | 559](#)
- [Verifying the DCBX Application Configuration | 560](#)
- [Verifying the DCBX Application Map Configuration | 560](#)
- [Verifying the DCBX Application Protocol Exchange Interface Configuration | 561](#)

To verify the configuration and proper operation of the lossless forwarding class and IEEE 802.1p priority, perform these tasks:

Verifying the Forwarding Class Configuration

Purpose

Verify that the lossless forwarding class `fcoe1` has been created.

Action

Show the forwarding class configuration by using the operational command `show class-of-service forwarding class`:

```
user@switch# show class-of-service forwarding-class
```

Forwarding class	ID	Queue	Policing priority	No-Loss
best-effort	0	0	normal	Disabled
fcoe	1	3	normal	Enabled
no-loss	2	4	normal	Enabled
network-control	3	7	normal	Disabled
fcoe1	4	5	normal	Enabled
mcast	8	8	normal	Disabled

Meaning

The `show class-of-service forwarding-class` command shows all of the forwarding classes. The command output shows that the `fcoe1` forwarding class is configured on output queue 5 with the no-loss packet drop attribute enabled.

Because we did not explicitly configure the default forwarding classes, they remain in their default state, including the lossless configuration of the `fcoe` and `no-loss` default forwarding classes.

Verifying the Behavior Aggregate Classifier Configuration

Purpose

Verify that the classifier maps the forwarding classes to the correct IEEE 802.1p code points (priorities) and packet loss priorities.

Action

List the classifier configured to support lossless FCoE transport using the operational mode command `show class-of-service classifier`:

```
user@switch> show class-of-service classifier
```

Classifier: fcoe_p5, Code point type: ieee-802.1, Index: 63065

Code point	Forwarding class	Loss priority
101	fcoe1	low

Meaning

The `show class-of-service classifier` command shows the IEEE 802.1p code points and the loss priorities that are mapped to the forwarding classes in each classifier.

Classifier `fcoe_p5` maps code point 101 (priority 5) to explicitly configured lossless forwarding class `fcoe1` and a packet loss priority of `low`, and all other priorities to the best-effort forwarding class with a packet loss priority of `high`.

Verifying the PFC Flow Control Configuration (CNP)

Purpose

Verify that PFC is enabled on the correct input priority and that flow control is configured on the correct output queue in the CNP.

Action

Display the congestion notification profile using the operational mode command `show class-of-service congestion-notification`:

```

user@switch> show class-of-service congestion-notification
Name: fcoe_p5_cnp, Index: 12137
Type: Input
Cable Length: 100 m
  Priority    PFC      MRU
  000        Disabled
  001        Disabled
  010        Disabled
  011        Disabled
  100        Disabled
  101        Enabled    2240
  110        Disabled
  111        Disabled
Type: Output
  Priority    Flow-Control-Queues
  101
                5

```

Meaning

The `show class-of-service congestion-notification` command shows the input and output stanzas of the configured CNPs.

The `fcoe_p5_cnp` CNP input stanza shows that PFC is enabled on code point 101 (priority 5), the MRU is 2240 bytes, and the cable length is 100 meters. The CNP output stanza shows that output flow control is configured on queue 5 for code point 101 (priority 5).

Verifying the Interface Configuration

Purpose

Verify that the correct classifier and congestion notification profile are configured on the interfaces.

Action

List the ingress interfaces using the operational mode commands `show configuration class-of-service interfaces xe-0/0/25` and `show configuration class-of-service interfaces xe-0/0/26`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/25
congestion-notification-profile fcoe_p5_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_p5;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/26
congestion-notification-profile fcoe_p5_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_p5;
    }
}
```

Meaning

Both the `show configuration class-of-service interfaces xe-0/0/25` command and the `show configuration class-of-service interfaces xe-0/0/26` command show that the congestion notification profile `fcoe_p5_cnp` is

configured on each interface, and that the IEEE 802.1p classifier associated with each interface is fcoe_p5.

Verifying the DCBX Application Configuration

Purpose

Verify that the DCBX application for FCoE is configured.

Action

List the DCBX applications by using the configuration mode command `show applications`:

```
user@switch# show applications
application fcoe_p5_app {
    ether-type 0x8906;
```

Meaning

The `show applications` configuration mode command shows all of the configured applications. The output shows that the application `fcoe_p5_app` is configured with an EtherType of `0x8906`.

Verifying the DCBX Application Map Configuration

Purpose

Verify that the application map is configured.

Action

List the application maps by using the configuration mode command `show policy-options application-maps`:

```
user@switch# show policy-options application-maps
fcoe_p5_app_map {
    application fcoe_p5_app code-points 101;
}
```

Meaning

The `show policy-options application-maps` configuration mode command lists all of the configured application maps and the applications that belong to each application map. The output shows that application map `fcoe_p5_app_map` consists of the application named `fcoe_p5_app`, which is mapped to IEEE 802.1p code point 101.

Verifying the DCBX Application Protocol Exchange Interface Configuration

Purpose

Verify that the application map is applied to the correct interfaces.

Action

List the application maps on each interface using the configuration mode command `show protocols dcbx`:

```
user@switch# show protocols dcbx
interface xe-0/0/25.0 {
    application-map fcoe_p5_app_map;
}
interface xe-0/0/26.0 {
    application-map fcoe_p5_app_map;
}
```

Meaning

The `show protocols dcbx` configuration mode command lists the application map association with interfaces. The output shows that interfaces `xe-0/0/25.0` and `xe-0/0/26.0` use application map `fcoe_p5_app_map`.

RELATED DOCUMENTATION

Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces

Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface

Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI)

Example: Configuring DCBX Application Protocol TLV Exchange

Configuring CoS PFC (Congestion Notification Profiles)

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface

IN THIS SECTION

- [Requirements | 562](#)
- [Overview | 563](#)
- [Configuration | 566](#)
- [Verification | 569](#)

The default system configuration supports FCoE traffic on priority 3 (IEEE 802.1p code point 011). If the FCoE traffic on your converged Ethernet network uses priority 3, the only user configuration required for lossless transport is to enable PFC on code point 011 on the FCoE ingress interfaces.

However, if your converged Ethernet network uses more than one priority for FCoE traffic, you need to configure lossless transport for each FCoE priority. This example shows you how to configure lossless FCoE transport on a converged Ethernet network that uses both priority 3 (IEEE 802.1p code point 011) and priority 5 (IEEE 802.1p code point 101) for FCoE traffic.

Requirements

This example uses the following hardware and software components:

- One switch used as an FCoE transit switch
- Junos OS Release 12.3 or later for the QFX Series

Overview

IN THIS SECTION

- [Topology | 564](#)

Some network topologies support FCoE traffic on more than one IEEE 802.1p priority. For example, a converged Ethernet network might include two separate FCoE networks that use different priorities to identify traffic. Interfaces that carry traffic for both FCoE networks need to support lossless FCoE transport on both priorities.

Supporting lossless behavior for two FCoE traffic classes requires configuring:

- At least one lossless forwarding class for FCoE traffic (this example uses the default `fcoe` forwarding class as one of the lossless FCoE forwarding classes, so we need to explicitly configure only one FCoE forwarding class).
- A behavior aggregate (BA) classifier to map the FCoE forwarding classes to the appropriate IEEE 802.1p code points (priorities).
- A congestion notification profile (CNP) to enable PFC on the FCoE code points at the interface ingress and to configure PFC flow control on the interface egress so that the interface can respond to PFC messages received from the connected peer.

NOTE: Configuring or changing PFC on an interface blocks the entire port until the PFC change is completed. After a PFC change is completed, the port is unblocked and traffic resumes. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

- DCBX applications and an application map to support DCBX application TLV exchange for the lossless FCoE traffic on the configured FCoE priorities. By default, DCBX is enabled on all Ethernet interfaces, but only on priority 3 (IEEE 802.1p code point 011). To support DCBX application TLV exchange when you are not using the default configuration, you must configure all of the applications and map them to interfaces and priorities.

The priorities specified in the BA classifier, CNP, and DCBX application map must match, or the configuration does not work. You must specify the same lossless FCoE forwarding class in each configuration and use the same IEEE 802.1p code point (priority) so that the FCoE traffic is properly classified into flows and so that those flows receive lossless treatment.

Topology

This example shows how to configure two lossless FCoE traffic classes on an interface, map them to two different priorities, and configure flow control to ensure lossless behavior. This example uses two Ethernet interfaces, xe-0/0/20 and xe-0/0/21, that are connected to the converged Ethernet network. Both interfaces transport FCoE traffic on priorities 3 (011) and 5 (101), and must support lossless transport of that traffic.

[Table 97 on page 564](#) shows the configuration components for this example.

Table 97: Components of the Two Lossless FCoE Priorities on an Interface Configuration Topology

Component	Settings
Hardware	One switch
Forwarding classes	<p>Name—fcoe1</p> <p>Queue mapping—queue 5</p> <p>Packet drop attribute—no-loss</p> <p>NOTE: A lossless forwarding class can be mapped to any output queue. However, because the fcoe1 forwarding class uses priority 5 in this example, matching that traffic to a forwarding class that uses queue 5 creates a configuration that is logical and easy to map because the priority and the queue are identified by the same number.</p> <p>Name—fcoe</p> <p>This is the default lossless FCoE forwarding class, so no configuration required. The fcoe forwarding class is mapped to priority 3 (IEEE 802.1p code point 011) and to output queue 3 with a packet drop attribute of no-loss.</p>
BA classifier	<p>Name—fcoe_classifier</p> <p>FCoE priority mapping for forwarding class fcoe—mapped to code point 011 (IEEE 802.1p priority 3) and a packet loss priority of low.</p> <p>FCoE priority mapping for forwarding class fcoe1—mapped to code point 101 (IEEE 802.1p priority 5) and a packet loss priority of low.</p>

Table 97: Components of the Two Lossless FCoE Priorities on an Interface Configuration Topology
(Continued)

Component	Settings
PFC configuration (CNP)	<p>CNP name—fcoe_cnp</p> <p>Input CNP code points—011 and 101</p> <p>MRU—2240 bytes</p> <p>Cable length—100 meters</p> <p>Output CNP code points—011 and 101</p> <p>Output CNP flow control queues—3 and 5</p> <p>NOTE: When you apply a CNP with an explicit output queue flow control configuration to an interface, the explicit CNP overwrites the default output CNP. The output queues that are enabled for PFC pause in the default configuration (queues 3 and 4) are not enabled for PFC pause unless they are included in the explicitly configured output CNP. In this example, because the explicit output CNP overwrites the default output CNP, we must explicitly configure flow control on queue 3.</p>
DCBX application mapping	<p>Application name—fcoe_app</p> <p>Application EtherType—0x8906</p> <p>Application map name—fcoe_app_map</p> <p>Application map code points—011 and 101</p> <p>NOTE: LLDP and DCBX must be enabled on the interface. By default, LLDP and DCBX are enabled on all Ethernet interfaces.</p>
Interfaces	<p>Interfaces xe-0/0/20 and xe-0/0/21 use the same configuration:</p> <ul style="list-style-type: none"> • Classifier—fcoe_classifier • CNP—fcoe_cnp • DCBX application map—fcoe_app_map

NOTE: This example does not include scheduling (bandwidth allocation) configuration or the FIP snooping configuration. This examples focuses only on the lossless FCoE priority configuration. QFX10000 switches do not support FIP snooping. For this reason, QFX10000 switches cannot be used as FCoE access transit switches. QFX10000 switches can be used as intermediate or aggregation transit switches in the FCoE path, between an FCoE access transit switch that performs FIP snooping and an FCF.

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 566](#)
- [Procedure | 567](#)

CLI Quick Configuration

To quickly configure two lossless FCoE forwarding classes that use different priorities on an FCoE transit switch interface, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level.

```
set class-of-service forwarding-classes class fcoe1 queue-num 5 no-loss
set class-of-service classifiers ieee-802.1 fcoe_classifier forwarding-class fcoe loss-priority
low code-points 011
set class-of-service classifiers ieee-802.1 fcoe_classifier forwarding-class fcoe1 loss-priority
low code-points 101set class-of-service interfaces xe-0/0/20 unit 0 classifiers ieee-802.1
fcoe_classifier
set class-of-service interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 fcoe_classifier
set class-of-service congestion-notification-profile fcoe_cnp input ieee-802.1 code-point 011
pfc mru 2240
set class-of-service congestion-notification-profile fcoe_cnp input ieee-802.1 code-point 101
pfc mru 2240
set class-of-service congestion-notification-profile fcoe_cnp input cable-length 100
set class-of-service congestion-notification-profile fcoe_cnp output ieee-802.1 code-point 011
pfc flow-control-queue 3
set class-of-service congestion-notification-profile fcoe_cnp output ieee-802.1 code-point 101
```

```
pfc flow-control-queue 5
set class-of-service interfaces xe-0/0/20 congestion-notification-profile fcoe_cnp
set class-of-service interfaces xe-0/0/21 congestion-notification-profile fcoe_cnp
set applications application fcoe_app ether-type 0x8906
set policy-options application-maps fcoe_app_map application fcoe_app code-points [011 101]
set protocols dcbx interface xe-0/0/20 application-map fcoe_app_map
set protocols dcbx interface xe-0/0/21 application-map fcoe_app_map
```

Procedure

Step-by-Step Procedure

To configure two lossless forwarding classes for FCoE traffic on the same interface, classify FCoE traffic into the forwarding classes, configure CNPs to enable PFC on the FCoE priorities and output queues, and configure DCBX application protocol TLV exchange for traffic on both FCoE priorities:

1. Configure lossless forwarding class `fcoe1` and map it to output queue 5 for FCoE traffic that uses IEEE 802.1p priority 5:

```
[edit class-of-service]
user@switch# set forwarding-classes class fcoe1 queue-num 5 no-loss
```

NOTE: This examples uses the default `fcoe` forwarding class as the other lossless FCoE forwarding class.

2. Configure the ingress classifier. The classifier maps the FCoE priorities (IEEE 802.1p code points 011 and 101) to lossless FCoE forwarding classes `fcoe` and `fcoe1`, respectively:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_classifier forwarding-class fcoe loss-priority low code-
points 011
user@switch# set ieee-802.1 fcoe_classifier forwarding-class fcoe1 loss-priority low code-
points 101
```

3. Apply the classifier to the interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 fcoe_classifier
user@switch# set interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 fcoe_classifier
```

4. Configure the CNP. The input stanza enables PFC on the FCoE priorities (IEEE 802.1p code points 011 and 101), sets the MRU value (2240 bytes), and sets the cable length value (100 meters). The output stanza configures flow control on output queues 3 and 5 on the FCoE priorities:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe_cnp input ieee-802.1 code-point 011 pfc
mru 2240
user@switch# set congestion-notification-profile fcoe_cnp input ieee-802.1 code-point 101 pfc
mru 2240
user@switch# set congestion-notification-profile fcoe_cnp input cable-length 100
user@switch# set congestion-notification-profile fcoe_cnp output ieee-802.1 code-point 011
pfc flow-control-queue 3
user@switch# set congestion-notification-profile fcoe_cnp output ieee-802.1 code-point 101
pfc flow-control-queue 5
```

5. Apply the CNP to the interfaces:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 congestion-notification-profile fcoe_cnp
user@switch# set interfaces xe-0/0/21 congestion-notification-profile fcoe_cnp
```

6. Configure a DCBX application for FCoE to map to the Ethernet interfaces, so that DCBX can exchange application protocol TLVs on both of the IEEE 802.1p priorities used for FCoE transport:

```
[edit]
user@switch# set applications application fcoe_app ether-type 0x8906
```

7. Configure a DCBX application map to map the FCoE application to the correct IEEE 802.1p FCoE priorities:

```
[edit]
user@switch# set policy-options application-maps fcoe_app_map application fcoe_app code-
points [011 101]
```

8. Apply the application map to the interfaces so that DCBX exchanges FCoE application TLVs on the correct code points:

```
[edit]
user@switch# set protocols dcbx interface xe-0/0/20 application-map fcoe_app_map
user@switch# set protocols dcbx interface xe-0/0/21 application-map fcoe_app_map
```

Verification

IN THIS SECTION

- [Verifying the Forwarding Class Configuration | 569](#)
- [Verifying the Behavior Aggregate Classifier Configuration | 570](#)
- [Verifying the PFC Flow Control Configuration \(CNP\) | 571](#)
- [Verifying the Interface Configuration | 572](#)
- [Verifying the DCBX Application Configuration | 573](#)
- [Verifying the DCBX Application Map Configuration | 573](#)
- [Verifying the DCBX Application Protocol Exchange Interface Configuration | 574](#)

To verify the configuration and proper operation of the lossless forwarding classes and IEEE 802.1p priorities, perform these tasks:

Verifying the Forwarding Class Configuration

Purpose

Verify that the lossless forwarding class `fcoe1` has been created.

Action

Show the forwarding class configuration by using the operational command `show class-of-service forwarding class`:

```
user@switch# show class-of-service forwarding-class
```

Forwarding class	ID	Queue	Policing priority	No-Loss
best-effort	0	0	normal	Disabled
fcoe	1	3	normal	Enabled
no-loss	2	4	normal	Enabled
network-control	3	7	normal	Disabled
fcoe1	4	5	normal	Enabled
mcast	8	8	normal	Disabled

Meaning

The `show class-of-service forwarding-class` command shows all of the forwarding classes. The command output shows that the `fcoe1` forwarding class is configured on output queue 5 with the no-loss packet drop attribute enabled.

Because we did not explicitly configure the default forwarding classes, they remain in their default state, including the lossless configuration of the `fcoe` and `no-loss` default forwarding classes.

Verifying the Behavior Aggregate Classifier Configuration

Purpose

Verify that the three classifiers map the forwarding classes to the correct IEEE 802.1p code points (priorities) and packet loss priorities.

Action

List the classifiers using the operational mode command `show class-of-service classifier`:

```
user@switch> show class-of-service classifier
```

Classifier: fcoe_classifier, Code point type: ieee-802.1, Index: 10964

Code point	Forwarding class	Loss priority
011	fcoe	low
101	fcoe1	low

Meaning

The `show class-of-service classifier` command shows the IEEE 802.1p code points and the loss priorities that are mapped to the forwarding classes in each classifier.

Classifier `fcoe_classifier` maps code point 011 to default lossless forwarding class `fcoe` and a packet loss priority of `low`, and maps code point 101 to explicitly configured lossless forwarding class `fcoe1` and a packet loss priority of `low`.

Verifying the PFC Flow Control Configuration (CNP)

Purpose

Verify that PFC is enabled on the correct input priorities and that flow control is configured on the correct output queues and priorities.

Action

List the CNPs using the operational mode command `show class-of-service congestion-notification`:

```
user@switch> show class-of-service congestion-notification
Name: fcoe_cnp, Index: 46504
Type: Input
Cable Length: 100 m
  Priority    PFC      MRU
  000        Disabled
  001        Disabled
  010        Disabled
  011        Enabled    2240
  100        Disabled
  101        Enabled    2240
  110        Disabled
  111        Disabled
Type: Output
  Priority    Flow-Control-Queues
  011
      3
  101
      5
```

Meaning

The `show class-of-service congestion-notification` command shows the input and output stanzas of the CNP.

The CNP `fcoe_cnp` input stanza shows that PFC is enabled on code points 011 and 101, the MRU is 2240 bytes on both priorities, and the interface cable length is 100 meters. The CNP output stanza shows that output flow control is configured on queues 3 and 5 for code points 011 and 101, respectively.

Verifying the Interface Configuration

Purpose

Verify that the classifier and congestion notification profile are configured on the interfaces. Both interfaces should show the same configuration.

Action

List the ingress interfaces using the operational mode commands `show configuration class-of-service interfaces xe-0/0/20` and `show configuration class-of-service interfaces xe-0/0/21`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/20
congestion-notification-profile fcoe_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_classifier;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/21
congestion-notification-profile fcoe_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_classifier;
    }
}
```

Meaning

The `show configuration class-of-service interfaces xe-0/0/20` command shows that the congestion notification profile `fcoe_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_classifier`.

The `show configuration class-of-service interfaces xe-0/0/21` command shows that the congestion notification profile `fcoe_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_classifier`.

Verifying the DCBX Application Configuration

Purpose

Verify that the DCBX application for FCoE is configured.

Action

List the DCBX applications by using the configuration mode command `show applications`:

```
user@switch# show applications
application fcoe_app {
    ether-type 0x8906;
```

Meaning

The `show applications` configuration mode command shows all of the configured applications. The output shows that the application `fcoe_app` is configured with an EtherType of `0x8906`.

Verifying the DCBX Application Map Configuration

Purpose

Verify that the application map is configured.

Action

List the application maps by using the configuration mode command `show policy-options application-maps`:

```
user@switch# show policy-options application-maps
fcoe_app_map {
    application fcoe_app code-points [011 101];
}
```

Meaning

The `show policy-options application-maps` configuration mode command lists all of the configured application maps and the applications that belong to each application map. The output shows that application map `fcoe_app_map` consists of the application named `fcoe_app`, which is mapped to IEEE 802.1p code points 011 and 101 (priorities 3 and 5, respectively).

Verifying the DCBX Application Protocol Exchange Interface Configuration

Purpose

Verify that the application map is applied to the interfaces.

Action

List the application maps on each interface using the configuration mode command `show protocols dcbx`:

```
user@switch# show protocols dcbx
interface xe-0/0/20.0 {
    application-map fcoe_app_map;
}
interface xe-0/0/21.0 {
    application-map fcoe_app_map;
}
```

Meaning

The `show protocols dcbx` configuration mode command lists the application map association with interfaces. The output shows that interfaces `xe-0/0/20.0` and `xe-0/0/21.0` use application map `fcoe_app_map`.

RELATED DOCUMENTATION

Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces

Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch)

Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI)

Example: Configuring DCBX Application Protocol TLV Exchange

Configuring CoS PFC (Congestion Notification Profiles)

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces

IN THIS SECTION

- [Requirements | 575](#)
- [Overview | 576](#)
- [Configuration | 581](#)
- [Verification | 586](#)

Although the default configuration provides two lossless forwarding classes mapped to two different IEEE 802.1p priorities (code points), you can explicitly configure up to six lossless forwarding classes and map them to different priorities. You can support up to six different types of lossless traffic, and you can support the same type of traffic if it uses different priorities in different parts of your converged network.

This example shows you how to configure two lossless forwarding classes for FCoE traffic and map them to two different priorities on an FCoE transit switch.

Requirements

This example uses the following hardware and software components:

- One switch used as an FCoE transit switch

- Junos OS Release 12.3 or later for the QFX Series

Overview

IN THIS SECTION

- [Topology | 577](#)

Some network topologies support FCoE traffic on more than one IEEE 802.1p priority. For example, when the switch acts as a transit switch, it could be connected to two QFX3500 switches in FCoE-FC gateway mode. Each of the gateway switches could connect a set of FCoE clients to a different SAN, and each set of FCoE clients could use a different priority for FCoE traffic to avoid fate sharing and maintain separation of the two FCoE networks. In this case, you need to configure two forwarding classes for FCoE traffic, each mapped to a different output queue and a different priority.

Supporting lossless behavior for two FCoE traffic classes requires configuring:

- At least one lossless forwarding class for FCoE traffic (this example uses the default `fcoe` forwarding class as one of the two lossless FCoE forwarding classes, so we need to explicitly configure only one FCoE forwarding class)
- Behavior aggregate (BA) classifiers to map the FCoE forwarding classes to the appropriate IEEE 802.1p code points (priorities) on each interface
- Congestion notification profiles (CNPs) for each interface to enable PFC on the FCoE code points at the interface ingress and to configure PFC flow control on the interface egress so that the interface can respond to PFC messages received from the connected peer

NOTE: Configuring or changing PFC on an interface blocks the entire port until the PFC change is completed. After a PFC change is completed, the port is unblocked and traffic resumes. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

- DCBX applications and an application map to support DCBX application TLV exchange for the lossless FCoE traffic on the configured FCoE priorities. By default, DCBX is enabled on all Ethernet interfaces, but only on priority 3 (IEEE 802.1p code point 011). To support DCBX application TLV exchange when you are not using the default configuration, you must configure all of the applications and map them to interfaces and priorities.

The priorities specified in the BA classifiers, CNPs, and DCBX application map must match, or the configuration does not work. You must specify the same lossless FCoE forwarding class in each configuration and use the same IEEE 802.1p code point (priority) so that the FCoE traffic is properly classified into flows and so that those flows receive lossless treatment.

Topology

This example shows how to configure two lossless FCoE traffic classes, map them to two different priorities, and configure flow control to ensure lossless behavior for those priorities on the interfaces. This example uses three Ethernet interfaces, xe-0/0/20, xe-0/0/21, and xe-0/0/22:

- Interface xe-0/0/20 connects to an FCoE-FC gateway that connects to Fibre Channel (FC) SAN 1. FCoE traffic to and from FC SAN 1 uses the default `fcoe` forwarding class and the default mapping to priority 3 (IEEE 802.1p code point 011) and output queue 3.
- Interface xe-0/0/21 connects to another FCoE-FC gateway that connects to Fibre Channel (FC) SAN 2. FCoE traffic to and from FC SAN-2 uses an explicitly configured FCoE forwarding class that is mapped to priority 5 (code point 101) and output queue 5.
- Interface xe-0/0/22 connects to FCoE devices on the converged Ethernet network and handles traffic destined for FC SAN 1 and FC SAN 2. Interface xe-0/0/22 must properly handle lossless FCoE traffic of both priorities (both FCoE forwarding classes), including pausing the traffic on ingress or egress as required.

Figure 23 on page 577 shows the topology for this example, and Table 98 on page 578 shows the configuration components for this example.

Figure 23: Topology of the Two Lossless FCoE Priorities Example

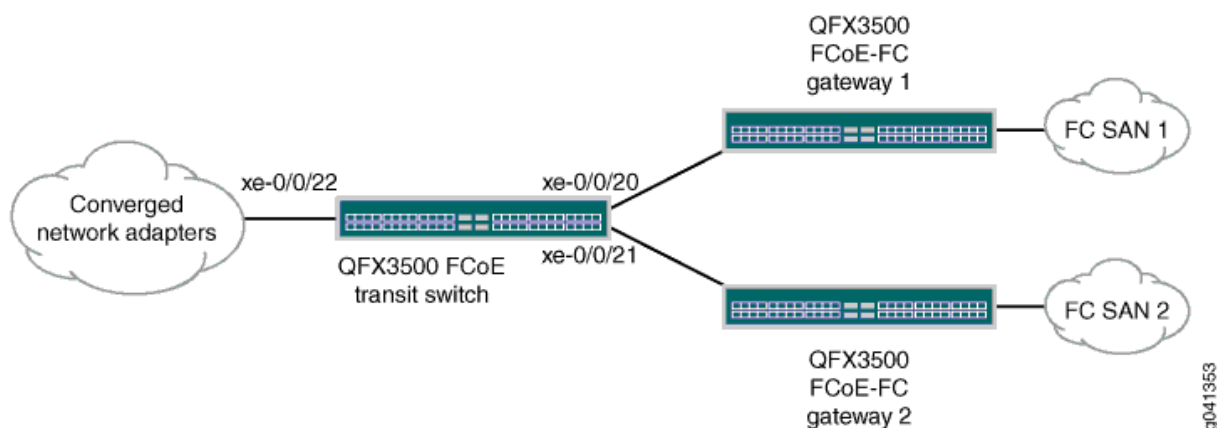


Table 98: Components of the Two Lossless FCoE Priorities Configuration Topology

Component	Settings
Hardware	One switch
Forwarding classes	<p>Name—fcoe1 Queue mapping—queue 5 Packet drop attribute—no-loss</p> <p>NOTE: A lossless forwarding class can be mapped to any output queue. However, because the fcoe1 forwarding class uses priority 5 in this example, matching that traffic to a forwarding class that uses queue 5 creates a configuration that is logical and easy to map because the priority and the queue are identified by the same number.</p> <p>Name—fcoe This is the default lossless FCoE forwarding class, so no configuration required. The fcoe forwarding class is mapped to priority 3 (IEEE 802.1p code point 011) and to output queue 3 with a packet drop attribute of no-loss</p>
BA classifiers	<p>Each interface requires a different classifier because each interface handles a different subset of FCoE traffic.</p> <ul style="list-style-type: none"> Interface xe-0/0/20 classifier: Name—fcoe_p3 FCoE priority mapping—Forwarding class fcoe mapped to code point 011 (IEEE 802.1p priority 3) and a packet loss priority of low. Interface xe-0/0/21 classifier: Name—fcoe_p5 FCoE priority mapping—Forwarding class fcoe1 mapped to code point 101 (IEEE 802.1p priority 5) and a packet loss priority of low. Interface xe-0/0/22 classifier: Name—fcoe_p3_p5 FCoE priority mapping—Forwarding class fcoe1 mapped to code point 101 and a packet loss priority of low, and forwarding class fcoe mapped to code point 011 and a packet loss priority of low.

Table 98: Components of the Two Lossless FCoE Priorities Configuration Topology (*Continued*)

Component	Settings
PFC configuration (CNPs)	<p>Each interface requires a different CNP because each interface handles a different subset of FCoE traffic and must pause that traffic on different priorities.</p> <ul style="list-style-type: none"> Interface xe-0/0/20 CNP: <ul style="list-style-type: none"> CNP name—fcoe_p3_cnp Input CNP code point—011 MRU—2240 bytes Cable length—100 meters <p>NOTE: Because interface xe-0/0/20 uses the default FCoE configuration, output queue 3 is paused by default and you do not need to configure the output stanza of the CNP.</p> Interface xe-0/0/21 CNP: <ul style="list-style-type: none"> CNP name—fcoe_p5_cnp Input CNP code point—101 MRU—2240 bytes Cable length—150 meters Output CNP code point—101 Output CNP flow control queue—5 Interface xe-0/0/22 CNP: <ul style="list-style-type: none"> CNP name—fcoe_p3_p5_cnp Input CNP code points—011 and 101 MRU—2240 bytes (both priorities) Cable length—100 meters Output CNP code points—011 (for queue 3) and 101 (for queue 5) Output CNP flow control queues—3 for priority 3 (code point 011) and 5 for priority 5 (code point 101) <p>NOTE: When you apply a CNP with an explicit output queue flow control configuration to an interface, the explicit CNP overwrites the default output CNP. The output queues that are enabled for pause in the default configuration (queues 3 and 4) are not enabled for pause unless they are included in the explicitly configured output CNP.</p>

Table 98: Components of the Two Lossless FCoE Priorities Configuration Topology (Continued)

Component	Settings
DCBX application mapping	<p>Interface xe-0/0/20 does not need an application map because DCBX exchanges application protocol TLVs only on the default FCoE priority (priority 3).</p> <p>Interface xe-0/0/21 requires an application map that enables DCBX application protocol TLV exchange on priority 5 (code point 101) for FCoE traffic. Interface xe-0/0/22 requires an application map that enables DCBX application protocol TLV exchange both on priority 3 (code point 011) and on priority 5 (code point 101) for FCoE traffic.</p> <ul style="list-style-type: none"> Interface xe-0/0/21 DCBX application mapping: <ul style="list-style-type: none"> Application name—fcoe_p5_app Application ether-type—0x8906 Application map name—fcoe_p5_app_map Application map code points—101 Interface xe-0/0/22 DCBX application mapping: <ul style="list-style-type: none"> Application name—fcoe_all_app Application ether-type—0x8906 Application map name—fcoe_all_app_map Application map code points—011 and 101 <p>NOTE: LLDP and DCBX must be enabled on the interface. By default, LLDP and DCBX are enabled on all Ethernet interfaces.</p>

NOTE: This example does not include scheduling (bandwidth allocation) configuration or the FIP snooping configuration. This examples focuses only on the lossless FCoE priority configuration. QFX10000 switches do not support FIP snooping. For this reason, QFX10000 switches cannot be used as FCoE access transit switches. QFX10000 switches can be used as intermediate or aggregation transit switches in the FCoE path, between an FCoE access transit switch that performs FIP snooping and an FCF.

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 581](#)
- [Procedure | 582](#)

CLI Quick Configuration

To quickly configure two lossless FCoE forwarding classes that use different priorities on an FCoE transit switch, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level.

```
set class-of-service forwarding-classes class fcoe1 queue-num 5 no-loss
set class-of-service classifiers ieee-802.1 fcoe_p3 forwarding-class fcoe loss-priority low code-
points 011
set class-of-service classifiers ieee-802.1 fcoe_p5 forwarding-class fcoe1 loss-priority low
code-points 101
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class fcoe loss-priority low
code-points 011
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class fcoe1 loss-priority low
code-points 101
set class-of-service interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 fcoe_p3
set class-of-service interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 fcoe_p5
set class-of-service interfaces xe-0/0/22 unit 0 classifiers ieee-802.1 fcoe_p3_p5
set class-of-service congestion-notification-profile fcoe_p3_cnp input ieee-802.1 code-point 011
pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p3_cnp input cable-length 100
set class-of-service congestion-notification-profile fcoe_p5_cnp input ieee-802.1 code-point 101
pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p5_cnp input cable-length 150
set class-of-service congestion-notification-profile fcoe_p5_cnp output ieee-802.1 code-point
101 pfc flow-control-queue 5
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp input ieee-802.1 code-point
011 pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp input ieee-802.1 code-point
101 pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp input cable-length 100
```



```

set class-of-service congestion-notification-profile fcoe_p3_p5_cnp output ieee-802.1 code-point
011 pfc flow-control-queue 3
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp output ieee-802.1 code-point
101 pfc flow-control-queue 5
set class-of-service interfaces xe-0/0/20 congestion-notification-profile fcoe_p3_cnp
set class-of-service interfaces xe-0/0/21 congestion-notification-profile fcoe_p5_cnp
set class-of-service interfaces xe-0/0/22 congestion-notification-profile fcoe_p3_p5_cnp
set applications application fcoe_p5_app ether-type 0x8906
set applications application fcoe_all_app ether-type 0x8906
set policy-options application-maps fcoe_p5_app_map application fcoe_p5_app code-points 101
set policy-options application-maps fcoe_all_app_map application fcoe_all_app code-points [011
101]
set protocols dcbx interface xe-0/0/21 application-map fcoe_p5_app_map
set protocols dcbx interface xe-0/0/22 application-map fcoe_all_app_map

```

Procedure

Step-by-Step Procedure

To configure two lossless forwarding classes for FCoE traffic on different interfaces, classify FCoE traffic into the forwarding classes, configure congestion notification profiles to enable PFC on the FCoE priorities and output queues, and configure DCBX application protocol TLV exchange for traffic on both FCoE priorities:

1. Configure lossless forwarding class fcoe1 and map it to output queue 5 for FCoE traffic that uses IEEE 802.1p priority 5:

```

[edit class-of-service]
user@switch# set forwarding-classes class fcoe1 queue-num 5 no-loss

```

NOTE: This examples uses the default fcoe forwarding class as the other lossless FCoE forwarding class.

2. Configure the ingress classifier (fcoe_p3) for interface xe-0/0/20. The classifier maps the FCoE priority (IEEE 802.1p code point 011) to lossless FCoE forwarding class fcoe:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_p3 forwarding-class fcoe loss-priority low code-points 011
```

3. Configure the ingress classifier (fcoe_p5) for interface xe-0/0/21. The classifier maps the FCoE priority (IEEE 802.1p code point 101) to lossless FCoE forwarding class fcoe1:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_p5 forwarding-class fcoe1 loss-priority low code-points 101
```

4. Configure the ingress classifier (fcoe_p3_p5) for interface xe-0/0/22. The classifier maps the two FCoE priorities (IEEE 802.1p code points 011 and 101) to the two lossless FCoE forwarding classes fcoe and fcoe1, respectively:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class fcoe loss-priority low code-points 011
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class fcoe1 loss-priority low code-points 101
```

5. Apply each classifier to the appropriate interface:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 fcoe_p3
user@switch# set interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 fcoe_p5
user@switch# set interfaces xe-0/0/22 unit 0 classifiers ieee-802.1 fcoe_p3_p5
```

6. Configure the CNP input stanza for interface xe-0/0/20 to enable PFC on the FCoE priority (IEEE 802.1p code point 011), set the MRU value (2240 bytes), and set the cable length value (100 meters). No output stanza is needed because queue 3 is paused by default on priority 3, and we are not explicitly configuring output queue flow control for any other queues.

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe_p3_cnp input ieee-802.1 code-point
```

```
011 pfc mru 2240
```

```
user@switch# set congestion-notification-profile fcoe_p3_cnp input cable-length 100
```

7. Configure the CNP for interface xe-0/0/21. The input stanza enables PFC on the FCoE priority (IEEE 802.1p code point 101), sets the MRU value (2240 bytes), and sets the cable length value (150 meters). The output stanza configures flow control on output queue 5 on the FCoE priority:

```
[edit class-of-service]
```

```
user@switch# set congestion-notification-profile fcoe_p5_cnp input ieee-802.1 code-point 101 pfc mru 2240
```

```
user@switch# set congestion-notification-profile fcoe_p5_cnp input cable-length 150
```

```
user@switch# set congestion-notification-profile fcoe_p5_cnp output ieee-802.1 code-point 101 pfc flow-control-queue 5
```

8. Configure the CNP for interface xe-0/0/22. The input stanza enables PFC on the FCoE priorities (IEEE 802.1p code points 011 and 101), sets the MRU value (2240 bytes), and sets the cable length value (100 meters). The output stanza configures flow control on output queues 3 and 5 on the FCoE priorities:

```
[edit class-of-service]
```

```
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp input ieee-802.1 code-point 011 pfc mru 2240
```

```
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp input ieee-802.1 code-point 101 pfc mru 2240
```

```
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp input cable-length 100
```

```
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp output ieee-802.1 code-point 011 pfc flow-control-queue 3
```

```
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp output ieee-802.1 code-point 101 pfc flow-control-queue 5
```

9. Apply each CNP to the appropriate interface:

```
[edit class-of-service]
```

```
user@switch# set interfaces xe-0/0/20 congestion-notification-profile fcoe_p3_cnp
```

```
user@switch# set interfaces xe-0/0/21 congestion-notification-profile fcoe_p5_cnp
```

```
user@switch# set interfaces xe-0/0/22 congestion-notification-profile fcoe_p3_p5_cnp
```

10. Configure the DCBX FCoE application and application map to apply to interface xe-0/0/21. Interface xe-0/0/21 uses priority 5 (IEEE 802.1p code point 101) for FCoE traffic, which requires DCBX to exchange FCoE application protocol TLVs on priority 5 on interface xe-0/0/21. Configure an application named `fcoe_p5_app` for FCoE traffic (EtherType 0x8906) and configure an application map named `fcoe_p5_app_map` to map the application to code point 101:

```
[edit]
user@switch# set applications application fcoe_p5_app ether-type 0x8906
user@switch# set policy-options application-maps fcoe_p5_app_map application fcoe_p5_app
code-points 101
```

NOTE: Interface xe-0/0/20 uses the default FCoE configuration (priority 3). DCBX exchanges protocol TLVs for the FCoE application by default, so you do not need to configure DCBX explicitly on interface xe-0/0/20.

11. Configure the DCBX FCoE application and application map to apply to interface xe-0/0/22. Interface xe-0/0/22 uses both priority 3 (IEEE 802.1p code point 011) and priority 5 for FCoE traffic, which requires DCBX to exchange FCoE application protocol TLVs on both priority 3 and priority 5. Configure an application named `fcoe_all_app` for FCoE traffic (EtherType 0x8906) and configure an application map named `fcoe_all_app_map` to map the application to code points 011 and 101:

```
[edit]
user@switch# set applications application fcoe_all_app ether-type 0x8906
user@switch# set policy-options application-maps fcoe_all_app_map application fcoe_all_app
code-points [011 101]
```

12. Apply the application maps to the interfaces xe-0/0/21 and xe-0/0/22 so that DCBX exchanges FCoE application TLVs on the correct code points on each interface:

```
[edit]
user@switch# set protocols dcbx interface xe-0/0/21 application-map fcoe_p5_app_map
user@switch# set protocols dcbx interface xe-0/0/22 application-map fcoe_all_app_map
```

Verification

IN THIS SECTION

- [Verifying the Forwarding Class Configuration | 586](#)
- [Verifying the Behavior Aggregate Classifier Configuration | 587](#)
- [Verifying the PFC Flow Control Configuration \(CNP\) | 588](#)
- [Verifying the Interface Configuration | 590](#)
- [Verifying the DCBX Application Configuration | 592](#)
- [Verifying the DCBX Application Map Configuration | 592](#)
- [Verifying the DCBX Application Protocol Exchange Interface Configuration | 593](#)

To verify the configuration and proper operation of the lossless forwarding classes and IEEE 802.1p priorities, perform these tasks:

Verifying the Forwarding Class Configuration

Purpose

Verify that the lossless forwarding class fcoe1 has been created.

Action

Show the forwarding class configuration by using the operational command `show class-of-service forwarding-class`:

```
user@switch# show class-of-service forwarding-class
```

Forwarding class	ID	Queue	Policing priority	No-Loss
best-effort	0	0	normal	Disabled
fcoe	1	3	normal	Enabled
no-loss	2	4	normal	Enabled
network-control	3	7	normal	Disabled
fcoe1	4	5	normal	Enabled
mcast	8	8	normal	Disabled

Meaning

The `show class-of-service forwarding-class` command shows all of the forwarding classes. The command output shows that the `fcoe1` forwarding class is configured on output queue 5 with the `no-loss` packet drop attribute enabled.

Because we did not explicitly configure the default forwarding classes, they remain in their default state, including the lossless configuration of the `fcoe` and `no-loss` default forwarding classes.

Verifying the Behavior Aggregate Classifier Configuration

Purpose

Verify that the three classifiers map the forwarding classes to the correct IEEE 802.1p code points (priorities) and packet loss priorities.

Action

List the classifiers configured to support lossless FCoE transport using the operational mode command `show class-of-service classifier`:

```
user@switch> show class-of-service classifier
Classifier: fcoe_p3, Code point type: ieee-802.1, Index: 13913
  Code point      Forwarding class      Loss priority
  ---
  011             fcoe                      low

Classifier: fcoe_p5, Code point type: ieee-802.1, Index: 63065
  Code point      Forwarding class      Loss priority
  ---
  101             fcoe1                     low

Classifier: fcoe_p3_p5, Code point type: ieee-802.1, Index: 10964
  Code point      Forwarding class      Loss priority
  ---
  011             fcoe                      low
  101             fcoe1                     low
```

Meaning

The `show class-of-service classifier` command shows the IEEE 802.1p code points and the loss priorities that are mapped to the forwarding classes in each classifier. The command output shows that there are three classifiers, `fcoe_p3`, `fcoe_p5`, and `fcoe_p3_p5`.

Classifier `fcoe_p3` maps code point 011 (priority 3) to default lossless forwarding class `fcoe` and a packet loss priority of `low`, and all other priorities to the best-effort forwarding class with a packet loss priority of `high`.

Classifier `fcoe_p5` maps code point 101 (priority 5) to explicitly configured lossless forwarding class `fcoe1` and a packet loss priority of `low`, and all other priorities to the best-effort forwarding class with a packet loss priority of `high`.

Classifier `fcoe_p3_p5` maps code point 011 to default lossless forwarding class `fcoe` and a packet loss priority of `low`, and maps code point 101 to explicitly configured lossless forwarding class `fcoe1` and a packet loss priority of `low`. The classifier maps all other priorities to the best-effort forwarding class with a packet loss priority of `high`.

Verifying the PFC Flow Control Configuration (CNP)

Purpose

Verify that PFC is enabled on the correct input priorities and that flow control is configured on the correct output queues and priorities in each CNP.

Action

List the congestion notification profiles using the operational mode command `show class-of-service congestion-notification`:

```
user@switch> show class-of-service congestion-notification
Name: fcoe_p3_cnp, Index: 12037
Type: Input
Cable Length: 100 m
  Priority    PFC          MRU
  000        Disabled
  001        Disabled
  010        Disabled
  011        Enabled    2240
  100        Disabled
  101        Disabled
  110        Disabled
  111        Disabled
Type: Output
  Priority    Flow-Control-Queues
  000
      0
```

001	1
010	2
011	3
100	4
101	5
110	6
111	7

Name: fcoe_p3_p5_cnp, Index: 46484

Type: Input

Cable Length: 100 m

Priority	PFC	MRU
000	Disabled	
001	Disabled	
010	Disabled	
011	Enabled	2240
100	Disabled	
101	Enabled	2240
110	Disabled	
111	Disabled	

Type: Output

Priority	Flow-Control-Queues
011	3
101	5

Name: fcoe_p5_cnp, Index: 12133

Type: Input

Cable Length: 150 m

Priority	PFC	MRU
000	Disabled	
001	Disabled	
010	Disabled	
011	Disabled	
100	Disabled	


```

101      Enabled      2240
110      Disabled
111      Disabled
Type: Output
Priority  Flow-Control-Queues
101
          5

```

Meaning

The `show class-of-service congestion-notification` command shows the input and output stanzas of the three CNPs. For CNP `fcoe_p3_cnp`, the input stanza shows that PFC is enabled on IEEE 802.1p code point 011 (priority 3), the MRU is 2240 bytes, and the cable length is 100 meters. The CNP output stanza shows the default mapping of priorities to output queues.

NOTE: By default, only queues 3 and 4 are enabled to respond to pause messages from the connected peer. For queue 3 to respond to pause messages, priority 3 (code point 011) must be enabled for PFC in the input stanza. For queue 4 to respond to pause messages, priority 4 (code point 100) must be enabled for PFC in the input stanza. In this example, only queue 3 responds to pause messages from the connected peer on interfaces that use CNP `fcoe_p3_cnp`, because the input stanza enables PFC priority 3 only.

For CNP `fcoe_p3_p5_cnp`, the input stanza shows that PFC is enabled on code points 011 and 101, the MRU is 2240 bytes on both priorities, and the cable length is 100 meters. The CNP output stanza shows that output flow control is configured on queues 3 and 5 for code points 011 and 101, respectively.

For CNP `fcoe_p5_cnp`, the input stanza shows that PFC is enabled on code point 101 (priority 5), the MRU is 2240 bytes, and the cable length is 150 meters. The CNP output stanza shows that output flow control is configured on queue 5 for code point 101 (priority 5).

Verifying the Interface Configuration

Purpose

Verify that the correct classifiers and congestion notification profiles are configured on the correct interfaces.

Action

List the ingress interfaces using the operational mode commands `show configuration class-of-service interfaces xe-0/0/20`, `show configuration class-of-service interfaces xe-0/0/21`, and `show configuration class-of-service interfaces xe-0/0/22`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/20
ccongestion-notification-profile fcoe_p3_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_p3;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/21
congestion-notification-profile fcoe_p5_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_p5;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/22
congestion-notification-profile fcoe_p3_p5_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_p3_p5;
    }
}
```

Meaning

The `show configuration class-of-service interfaces xe-0/0/20` command shows that the congestion notification profile `fcoe_p3_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_p3`.

The `show configuration class-of-service interfaces xe-0/0/21` command shows that the congestion notification profile `fcoe_p5_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_p5`.

The `show configuration class-of-service interfaces xe-0/0/22` command shows that the congestion notification profile `fcoe_p3_p5_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_p3_p5`.

Verifying the DCBX Application Configuration

Purpose

Verify that the two DCBX applications for FCoE are configured.

Action

List the DCBX applications by using the configuration mode command `show applications`:

```
user@switch# show applications
application fcoe_all_app {
    ether-type 0x8906;

application fcoe_p5_app {
    ether-type 0x8906;
```

Meaning

The `show applications` configuration mode command shows all of the configured applications. The output shows that the application `fcoe_all_app` is configured with an EtherType of `0x8906` (the correct EtherType for FCoE traffic) and that the application `fcoe_p5_app` is also configured with an EtherType of `0x8906`.

Verifying the DCBX Application Map Configuration

Purpose

Verify that the application maps are configured.

Action

List the application maps by using the configuration mode command `show policy-options application-maps`:

```
user@switch# show policy-options application-maps
fcoe_all_app_map {
```

```

    application fcoe_all_app code-points [011 101];
}
fcoe_p5_app_map {
    application fcoe_p5_app code-points 101;
}

```

Meaning

The `show policy-options application-maps` configuration mode command lists all of the configured application maps and the applications that belong to each application map. The output shows that there are two application maps.

Application map `fcoe_all_app_map` consists of the application named `fcoe_all_app` mapped to IEEE 802.1p code points 011 (priority 3) and 101 (priority 5).

Application map `fcoe_p5_app_map` consists of the application named `fcoe_p5_app` mapped to IEEE 802.1p code point 101 (priority 5).

Verifying the DCBX Application Protocol Exchange Interface Configuration

Purpose

Verify that the application maps are applied to the correct interfaces.

Action

List the application maps on each interface using the configuration mode command `show protocols dcbx`:

```

user@switch# show protocols dcbx
interface xe-0/0/21.0 {
    application-map fcoe_p5_app_map;
}
interface xe-0/0/22.0 {
    application-map fcoe_all_app_map;
}

```

Meaning

The `show protocols dcbx` configuration mode command lists the application map association with interfaces. The output shows that interface `xe-0/0/21.0` uses application map `fcoe_p5_app_map` and interface `xe-0/0/22.0` uses application map `fcoe_all_app_map`.

NOTE: Because interface xe-0/0/20 uses the default lossless FCoE configuration, you do not configure application mapping to interface xe-0/0/20. The default configuration automatically exchanges application protocol TLVs for the default FCoE configuration on priority 3 (IEEE 802.1p code point 011).

RELATED DOCUMENTATION

Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface

Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch)

Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI)

Example: Configuring DCBX Application Protocol TLV Exchange

Configuring CoS PFC (Congestion Notification Profiles)

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI)

IN THIS SECTION

- [Requirements | 595](#)
- [Overview | 595](#)
- [Configuration | 602](#)
- [Verification | 610](#)

Although the default configuration provides two lossless forwarding classes mapped to two different IEEE 802.1p priorities (code points), you can explicitly configure up to six lossless forwarding classes and

map them to different priorities. You can support up to six different types of lossless traffic, and you can support the same type of traffic on different priorities in different parts of your converged network.

This example shows you how to configure two lossless forwarding classes for FCoE traffic and one lossless forwarding class for iSCSI traffic, and map the forwarding classes to three different priorities. (The converged Ethernet network includes two FCoE networks, each of which uses a different priority to identify FCoE traffic, and an iSCSI network.)

Requirements

This example uses the following hardware and software components:

- One switch used as an FCoE transit switch
- Junos OS Release 12.3 or later for the QFX Series

Overview

IN THIS SECTION

- [Topology | 596](#)

Some converged Ethernet networks support FCoE on more than one IEEE 802.1p priority and also require supporting other lossless traffic classes. Interfaces that carry multiple lossless forwarding classes need to support lossless behavior for the priorities mapped to those forwarding classes. To support the two FCoE forwarding classes and the iSCSI forwarding class used in this example, you need to configure:

- At least one lossless forwarding class for FCoE traffic (this example uses the default `fcoe` forwarding class as one of the two lossless FCoE forwarding classes, so we need to explicitly configure only one FCoE forwarding class)
- A lossless forwarding class for iSCSI traffic
- Behavior aggregate (BA) classifiers to map the lossless forwarding classes to the appropriate IEEE 802.1p code points (priorities) on each interface
- Congestion notification profiles (CNPs) for each interface to enable PFC on the FCoE and iSCSI code points at the interface ingress, and to configure PFC flow control on the interface egress so that the interface can respond to PFC messages received from the connected peer

NOTE: Configuring or changing PFC on an interface blocks the entire port until the PFC change is completed. After a PFC change is completed, the port is unblocked and traffic resumes. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

- DCBX applications and an application map to support DCBX application TLV exchange for the FCoE and iSCSI traffic on the configured lossless priorities. By default, DCBX is enabled on all Ethernet interfaces for FCoE, but only on priority 3 (IEEE 802.1p code point 011). To support DCBX application TLV exchange when you are not using the default configuration, you must configure all of the applications and map them to interfaces and priorities.

The priorities specified in the BA classifiers, CNPs, and DCBX application map must match, or the configuration does not work. You must specify the same lossless FCoE forwarding class in each configuration and use the same IEEE 802.1p code point (priority) so that the FCoE traffic is properly classified into flows and so that those flows receive lossless treatment.

Topology

This example shows how to configure two lossless FCoE traffic classes and one lossless iSCSI traffic class, map them to three different priorities, and configure flow control to ensure lossless behavior for those priorities on the interfaces. This example uses four Ethernet interfaces, xe-0/0/31, xe-0/0/32, xe-0/0/33, and xe-0/0/34:

- Interface xe-0/0/31 handles FCoE traffic on priority 3 (IEEE 802.1p code point 011) and iSCSI traffic on priority 4 (code point 100).
- Interface xe-0/0/32 handles FCoE traffic on priority 5 (code point 101) and iSCSI traffic on priority 4.
- Interface xe-0/0/33 handles FCoE traffic on priority 3 and priority 5.
- Interface xe-0/0/34 handles iSCSI traffic on priority 4.

[Figure 24 on page 597](#) shows the topology for this example, and [Table 99 on page 597](#) shows the configuration components for this example.

Figure 24: Topology of the Lossless FCoE and iSCSI Priorities Example

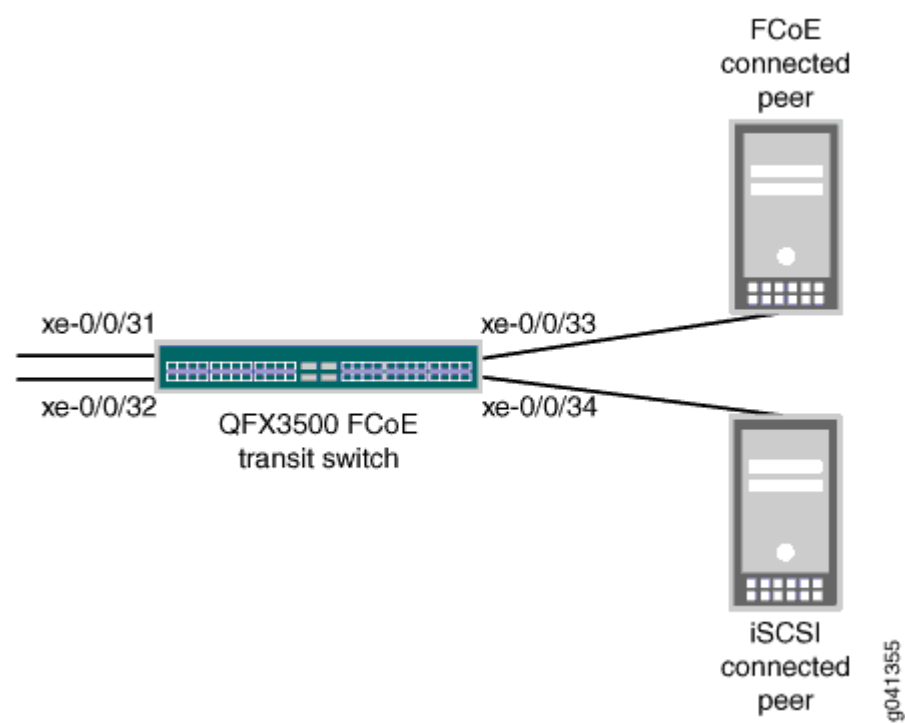


Table 99: Components of the Lossless FCoE and iSCSI Priorities Configuration Topology

Component	Settings
Hardware	One switch

Table 99: Components of the Lossless FCoE and iSCSI Priorities Configuration Topology (*Continued*)

Component	Settings
Forwarding classes	<p>This example uses one explicitly configured lossless FCoE forwarding class, the default lossless FCoE forwarding class, and one explicitly configured iSCSI forwarding class.</p> <ul style="list-style-type: none"> iSCSI forwarding class: <ul style="list-style-type: none"> Name—iscsi Queue mapping—queue 4 Packet drop attribute—no-loss FCoE forwarding class (explicitly configured): <ul style="list-style-type: none"> Name—fcoe1 Queue mapping—queue 5 Packet drop attribute—no-loss <p>NOTE: A lossless forwarding class can be mapped to any output queue. However, because the fcoe1 forwarding class uses priority 5 in this example, matching that traffic to a forwarding class that uses queue 5 creates a configuration that is logical and easy to map because the priority and the queue are identified by the same number.</p> <ul style="list-style-type: none"> FCoE forwarding class (default) <ul style="list-style-type: none"> Name—fcoe <p>The default fcoe forwarding class is mapped to priority 3 (IEEE 802.1p code point 011) and to output queue 3 with a packet drop attribute of no-loss.</p>

Table 99: Components of the Lossless FCoE and iSCSI Priorities Configuration Topology (*Continued*)

Component	Settings
BA classifiers	<p>Each interface requires a different classifier because each interface handles a different subset of FCoE traffic.</p> <ul style="list-style-type: none"> <p>Interface xe-0/0/31 classifier:</p> <p>Name—fcoe_p3_iscsi</p> <p>FCoE priority mapping—Forwarding class fcoe mapped to code point 011 (IEEE 802.1p priority 3) and a packet loss priority of low.</p> <p>iSCSI priority mapping—Forwarding class iscsi mapped to code point 100 (priority 4) and a packet loss priority of low.</p> <p>All other priority mapping—All other forwarding classes are mapped to the best-effort forwarding class with packet loss priorities of high.</p> <p>Interface xe-0/0/32 classifier:</p> <p>Name—fcoe_p5_iscsi</p> <p>FCoE priority mapping—Forwarding class fcoe1 mapped to code point 101 (IEEE 802.1p priority 5) and a packet loss priority of low.</p> <p>iSCSI priority mapping—Forwarding class iscsi mapped to code point 100 (priority 4) and a packet loss priority of low.</p> <p>All other priority mapping—All other forwarding classes are mapped to the best-effort forwarding class with packet loss priorities of high.</p> <p>Interface xe-0/0/33 classifier:</p> <p>Name—fcoe_p3_p5</p> <p>FCoE priority mapping—Forwarding class fcoe1 mapped to code point 101 (priority 5) and a packet loss priority of low, and forwarding class fcoe mapped to code point 011 and a packet loss priority of low.</p> <p>All other priority mapping—All other forwarding classes are mapped to the best-effort forwarding class with packet loss priorities of high.</p> <p>Interface xe-0/0/34 classifier:</p> <p>Name—iscsi_classifier</p> <p>iSCSI priority mapping—Forwarding class iscsi mapped to code point 100 (priority 4) and a packet loss priority of low.</p> <p>All other priority mapping—All other forwarding classes are mapped to the best-effort forwarding class with packet loss priorities of high.</p>

Table 99: Components of the Lossless FCoE and iSCSI Priorities Configuration Topology (*Continued*)

Component	Settings
PFC configuration (CNPs)	<p>Each interface requires a different CNP because each interface handles a different subset of FCoE and iSCSI traffic, and must pause that traffic on different priorities.</p> <ul style="list-style-type: none"> Interface xe-0/0/31 CNP: <ul style="list-style-type: none"> CNP name—fcoe_p3_cnp Input CNP code points—011 and 100 MRU—2240 bytes for code point 011, default value (2500 bytes) for code point 100 Cable length—100 meters <p>NOTE: On interface xe-0/0/31, the FCoE forwarding class is mapped to queue 3 and priority 3 (code point 011), and the iSCSI forwarding class is mapped to queue 4 and priority 4 (code point 100). Therefore, interface xe-0/0/31 does not require an output CNP configuration because queue 3 and queue 4 are enabled for PFC flow control by default on code points 011 and 100, respectively.</p> Interface xe-0/0/32 CNP: <ul style="list-style-type: none"> CNP name—fcoe_p5_cnp Input CNP code points—100 and 101 MRU—Default value (2500 bytes) for code point 100, 2240 bytes for code point 101 Cable length—150 meters Output CNP code points—100 and 101 Output CNP flow control queues—4 and 5 Interface xe-0/0/33 CNP: <ul style="list-style-type: none"> CNP name—fcoe_p3_p5_cnp Input CNP code points—011 and 101 MRU—2240 bytes (both priorities) Cable length—100 meters Output CNP code points—011 and 101 Output CNP flow control queues—3 and 5 Interface xe-0/0/34 CNP: <ul style="list-style-type: none"> CNP name—iscsi_cnp Input CNP code point—100 MRU—2500 bytes (default value)

Table 99: Components of the Lossless FCoE and iSCSI Priorities Configuration Topology (*Continued*)

Component	Settings
	<p>Cable length—100 meters</p> <p>NOTE: On interface xe-0/0/34, the iSCSI forwarding class is mapped to queue 4 and priority 4 (code point 100). Interface xe-0/0/34 does not require an output CNP configuration because queue 4 is enabled for PFC flow control by default on code point 100.</p> <p>NOTE: When you apply a CNP with an explicit output queue flow control configuration to an interface, the explicit CNP overwrites the default output CNP. The output queues that are enabled for PFC pause in the default configuration (queues 3 and 4) are not enabled for pause unless they are included in the explicitly configured output CNP.</p>
DCBX application mapping	<p>This example requires configuring applications for FCoE and iSCSI, including them in the same application map, and applying the application map to all four interfaces.</p> <p>Application map name—<code>dcbx_iscsi_fcoe_app_map</code></p> <ul style="list-style-type: none"> FCoE application name—<code>fcoe_app</code> Application ether-type—<code>0x8906</code> Application map code points—<code>011</code> and <code>101</code> iSCSI application name—<code>iscsi_app</code> Application protocol type—<code>tcp</code> Application destination port—<code>3260</code> Application map code point—<code>100</code> <p>NOTE: LLDP and DCBX must be enabled on the interface. By default, LLDP and DCBX are enabled on all Ethernet interfaces.</p>

NOTE: This example does not include scheduling (bandwidth allocation) configuration or the FIP snooping configuration. This examples focuses only on the lossless FCoE priority configuration. QFX10000 switches do not support FIP snooping. For this reason, QFX10000 switches cannot be used as FCoE access transit switches. QFX10000 switches can be used as intermediate or aggregation transit switches in the FCoE path, between an FCoE access transit switch that performs FIP snooping and an FCF.

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 602](#)
- [Procedure | 605](#)

CLI Quick Configuration

To quickly configure two lossless FCoE forwarding classes and one lossless iSCSI forwarding class and map them to different priorities, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level.

```
set class-of-service forwarding-classes class iscsi queue-num 4 no-loss
set class-of-service forwarding-classes class fcoe1 queue-num 5 no-loss
set class-of-service classifiers ieee-802.1 fcoe_p3_iscsi forwarding-class fcoe loss-priority
low code-points 011
set class-of-service classifiers ieee-802.1 fcoe_p3_iscsi forwarding-class iscsi loss-priority
low code-points 100
set class-of-service classifiers ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-
priority high code-points 000
set class-of-service classifiers ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-
priority high code-points 001
set class-of-service classifiers ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-
priority high code-points 010
set class-of-service classifiers ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-
priority high code-points 101
set class-of-service classifiers ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-
priority high code-points 110
set class-of-service classifiers ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-
priority high code-points 111
set class-of-service classifiers ieee-802.1 fcoe_p5_iscsi forwarding-class iscsi loss-priority
low code-points 100
set class-of-service classifiers ieee-802.1 fcoe_p5_iscsi forwarding-class fcoe1 loss-priority
low code-points 101
set class-of-service classifiers ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-
priority high code-points 000
set class-of-service classifiers ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-
```

```

priority high code-points 001
set class-of-service classifiers ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-
priority high code-points 010
set class-of-service classifiers ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-
priority high code-points 011
set class-of-service classifiers ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-
priority high code-points 110
set class-of-service classifiers ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-
priority high code-points 111
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class fcoe loss-priority low
code-points 011
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class fcoe1 loss-priority low
code-points 101
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-
priority high code-points 000
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-
priority high code-points 001
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-
priority high code-points 010
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-
priority high code-points 100
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-
priority high code-points 110
set class-of-service classifiers ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-
priority high code-points 111
set class-of-service classifiers ieee-802.1 iscsi_classifier forwarding-class iscsi loss-
priority low code-points 100
set class-of-service classifiers ieee-802.1 iscsi_classifier forwarding-class best-effort loss-
priority high code-points 000
set class-of-service classifiers ieee-802.1 iscsi_classifier forwarding-class best-effort loss-
priority high code-points 001
set class-of-service classifiers ieee-802.1 iscsi_classifier forwarding-class best-effort loss-
priority high code-points 010
set class-of-service classifiers ieee-802.1 iscsi_classifier forwarding-class best-effort loss-
priority high code-points 011
set class-of-service classifiers ieee-802.1 iscsi_classifier forwarding-class best-effort loss-
priority high code-points 101
set class-of-service classifiers ieee-802.1 iscsi_classifier forwarding-class best-effort loss-
priority high code-points 110
set class-of-service classifiers ieee-802.1 iscsi_classifier forwarding-class best-effort loss-
priority high code-points 111
set class-of-service interfaces xe-0/0/31 unit 0 classifiers ieee-802.1 fcoe_p3_iscsi
set class-of-service interfaces xe-0/0/32 unit 0 classifiers ieee-802.1 fcoe_p5_iscsi

```

```

set class-of-service interfaces xe-0/0/33 unit 0 classifiers ieee-802.1 fcoe_p3_p5
set class-of-service interfaces xe-0/0/34 unit 0 classifiers ieee-802.1 iscsi_classifier
set class-of-service congestion-notification-profile fcoe_p3_cnp input ieee-802.1 code-point 011
pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p3_cnp input ieee-802.1 code-point 100
pfc
set class-of-service congestion-notification-profile fcoe_p3_cnp input cable-length 100
set class-of-service congestion-notification-profile fcoe_p5_cnp input ieee-802.1 code-point 100
pfc
set class-of-service congestion-notification-profile fcoe_p5_cnp input ieee-802.1 code-point 101
pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p5_cnp input cable-length 150
set class-of-service congestion-notification-profile fcoe_p5_cnp output ieee-802.1 code-point
100 pfc flow-control-queue 4
set class-of-service congestion-notification-profile fcoe_p5_cnp output ieee-802.1 code-point
101 pfc flow-control-queue 5
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp input ieee-802.1 code-point
011 pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp input ieee-802.1 code-point
101 pfc mru 2240
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp input cable-length 100
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp output ieee-802.1 code-point
011 pfc flow-control-queue 3
set class-of-service congestion-notification-profile fcoe_p3_p5_cnp output ieee-802.1 code-point
101 pfc flow-control-queue 5
set class-of-service congestion-notification-profile iscsi_cnp input ieee-802.1 code-point 100
pfc
set class-of-service congestion-notification-profile iscsi_cnp input cable-length 100
set class-of-service interfaces xe-0/0/31 congestion-notification-profile fcoe_p3_cnp
set class-of-service interfaces xe-0/0/32 congestion-notification-profile fcoe_p5_cnp
set class-of-service interfaces xe-0/0/33 congestion-notification-profile fcoe_p3_p5_cnp
set class-of-service interfaces xe-0/0/34 congestion-notification-profile iscsi_cnp
set applications application iscsi_app protocol tcp destination-port 3260
set applications application fcoe_app ether-type 0x8906
set policy-options application-maps dcbx_iscsi_fcoe_app_map application iscsi_app code-points 100
set policy-options application-maps dcbx_iscsi_fcoe_app_map application fcoe_app code-points
[011 101]
set protocols dcbx interface xe-0/0/31 application-map dcbx_iscsi_fcoe_app_map
set protocols dcbx interface xe-0/0/32 application-map dcbx_iscsi_fcoe_app_map
set protocols dcbx interface xe-0/0/33 application-map dcbx_iscsi_fcoe_app_map
set protocols dcbx interface xe-0/0/34 application-map dcbx_iscsi_fcoe_app_map

```

Procedure

Step-by-Step Procedure

To configure two lossless forwarding classes for FCoE traffic and one lossless forwarding class for iSCSI traffic, classify the traffic into the three forwarding classes, configure congestion notification profiles to enable PFC on the FCoE priorities and output queues, and configure DCBX application protocol TLV exchange for traffic on both FCoE priorities:

1. Configure lossless forwarding classes `iscsi` for iSCSI traffic and `fcoe1` for FCoE traffic (this example uses the default `fcoe` forwarding class as the other lossless FCoE forwarding class) and map them to output queues:

```
[edit class-of-service]
user@switch# set forwarding-classes class iscsi queue-num 4 no-loss
user@switch# set forwarding-classes class fcoe1 queue-num 5 no-loss
```

2. Configure the ingress classifier (`fcoe_p3_iscsi`) for interface `xe-0/0/31`. The classifier maps the FCoE priority (code point 011) to lossless FCoE forwarding class `fcoe` and the iSCSI priority (code point 100) to lossless iSCSI forwarding class `iscsi`, and traffic of other priorities to the best-effort forwarding class with a packet loss priority of high:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_p3_iscsi forwarding-class fcoe loss-priority low code-points 011
user@switch# set ieee-802.1 fcoe_p3_iscsi forwarding-class iscsi loss-priority low code-points 100
user@switch# set ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-priority high code-points 000
user@switch# set ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-priority high code-points 001
user@switch# set ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-priority high code-points 010
user@switch# set ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-priority high code-points 101
user@switch# set ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-priority high code-points 110
user@switch# set ieee-802.1 fcoe_p3_iscsi forwarding-class best-effort loss-priority high code-points 111
```


3. Configure the ingress classifier (fcoe_p5_iscsi) for interface xe-0/0/32. The classifier maps the FCoE priority (code point 101) to lossless FCoE forwarding class fcoe1 and the iSCSI priority (code point 100) to lossless iSCSI forwarding class iscsi, and traffic of other priorities to the best-effort forwarding class with a packet loss priority of high:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_p5_iscsi forwarding-class iscsi loss-priority low code-
points 100
user@switch# set ieee-802.1 fcoe_p5_iscsi forwarding-class fcoe1 loss-priority low code-
points 101
user@switch# set ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-priority high
code-points 000
user@switch# set ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-priority high
code-points 001
user@switch# set ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-priority high
code-points 010
user@switch# set ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-priority high
code-points 011
user@switch# set ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-priority high
code-points 110
user@switch# set ieee-802.1 fcoe_p5_iscsi forwarding-class best-effort loss-priority high
code-points 111
```

4. Configure the ingress classifier (fcoe_p3_p5) for interface xe-0/0/33. The classifier maps the two FCoE priorities (code points 011 and 101) to lossless FCoE forwarding classes fcoe and fcoe1, respectively, and traffic of other priorities to the best-effort forwarding class with a packet loss priority of high:

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class fcoe loss-priority low code-points
011
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class fcoe1 loss-priority low code-points
101
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-priority high code-
points 000
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-priority high code-
points 001
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-priority high code-
points 010
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-priority high code-
points 100
```

```

user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-priority high code-
points 110
user@switch# set ieee-802.1 fcoe_p3_p5 forwarding-class best-effort loss-priority high code-
points 111

```

5. Configure the ingress classifier (iscsi_classifier) for interface xe-0/0/34. The classifier maps the iSCSI priority (code point 101) to lossless iSCSI forwarding class iscsi, and traffic of other priorities to the best-effort forwarding class with a packet loss priority of high:

```

[edit class-of-service classifiers]
user@switch# set ieee-802.1 iscsi_classifier forwarding-class iscsi loss-priority low code-
points 100
user@switch# set ieee-802.1 iscsi_classifier forwarding-class best-effort loss-priority
high code-points 000
user@switch# set ieee-802.1 iscsi_classifier forwarding-class best-effort loss-priority
high code-points 001
user@switch# set ieee-802.1 iscsi_classifier forwarding-class best-effort loss-priority
high code-points 010
user@switch# set ieee-802.1 iscsi_classifier forwarding-class best-effort loss-priority
high code-points 011
user@switch# set ieee-802.1 iscsi_classifier forwarding-class best-effort loss-priority
high code-points 101
user@switch# set ieee-802.1 iscsi_classifier forwarding-class best-effort loss-priority
high code-points 110
user@switch# set ieee-802.1 iscsi_classifier forwarding-class best-effort loss-priority
high code-points 111

```

6. Apply each classifier to the appropriate interface:

```

[edit class-of-service]
user@switch# set interfaces xe-0/0/31 unit 0 classifiers ieee-802.1 fcoe_p3_iscsi
user@switch# set interfaces xe-0/0/32 unit 0 classifiers ieee-802.1 fcoe_p5_iscsi
user@switch# set interfaces xe-0/0/33 unit 0 classifiers ieee-802.1 fcoe_p3_p5
user@switch# set interfaces xe-0/0/34 unit 0 classifiers ieee-802.1 iscsi_classifier

```

7. Configure the CNP input stanza for interface xe-0/0/31 to enable PFC on the FCoE and iSCSI priorities that the interface handles (code points 011 and 100), set the MRU value for the FCoE traffic (2240 bytes), and set the cable length value (100 meters). No output stanza is needed

because queues 3 and 4 are paused by default on priorities 3 and 4, respectively, and we are not explicitly configuring output queue flow control for any other queues.

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe_p3_cnp input ieee-802.1 code-point
011 pfc mru 2240
user@switch# set congestion-notification-profile fcoe_p3_cnp input ieee-802.1 code-point
100 pfc
user@switch# set congestion-notification-profile fcoe_p3_cnp input cable-length 100
```

8. Configure the CNP for interface xe-0/0/32. The input stanza enables PFC on the FCoE priority (code point 101), sets the MRU value for FCoE traffic (2240 bytes), enables PFC on the iSCSI priority (code point 100), and sets the cable length value (150 meters). The output stanza configures flow control on output queue 5 on the FCoE priority and on output queue 4 on the iSCSI priority:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe_p5_cnp input ieee-802.1 code-point
100 pfc
user@switch# set congestion-notification-profile fcoe_p5_cnp input ieee-802.1 code-point
101 pfc mru 2240
user@switch# set congestion-notification-profile fcoe_p5_cnp input cable-length 150
user@switch# set congestion-notification-profile fcoe_p5_cnp output ieee-802.1 code-point
100 pfc flow-control-queue 4
user@switch# set congestion-notification-profile fcoe_p5_cnp output ieee-802.1 code-point
101 pfc flow-control-queue 5
```

9. Configure the CNP for interface xe-0/0/33. The input stanza enables PFC on the FCoE priorities (IEEE 802.1p code points 011 and 101), sets the MRU value (2240 bytes), and sets the cable length value (100 meters). The output stanza configures flow control on output queues 3 and 5 on the FCoE priorities:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp input ieee-802.1 code-point
011 pfc mru 2240
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp input ieee-802.1 code-point
101 pfc mru 2240
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp input cable-length 100
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp output ieee-802.1 code-
point 011 pfc flow-control-queue 3
```

```
user@switch# set congestion-notification-profile fcoe_p3_p5_cnp output ieee-802.1 code-
point 101 pfc flow-control-queue 5
```

10. Configure the CNP input stanza for interface xe-0/0/34 to enable PFC on the iSCSI priority (code point 100) and set the cable length value (100 meters). No output stanza is needed because queue 4 is paused by default on priority 4, and we are not explicitly configuring output queue flow control for any other queues.

```
[edit class-of-service]
user@switch# set congestion-notification-profile iscsi_cnp input ieee-802.1 code-point 100
pfc
user@switch# set congestion-notification-profile iscsi_cnp input cable-length 100
```

11. Apply each CNP to the appropriate interface:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/31 congestion-notification-profile fcoe_p3_cnp
user@switch# set interfaces xe-0/0/32 congestion-notification-profile fcoe_p5_cnp
user@switch# set interfaces xe-0/0/33 congestion-notification-profile fcoe_p3_p5_cnp
user@switch# set interfaces xe-0/0/34 congestion-notification-profile iscsi_cnp
```

12. Configure the DCBX applications for FCoE and iSCSI to map to the interfaces so that DCBX can exchange application protocol TLVs on the IEEE 802.1p priorities used for FCoE and iSCSI traffic:

```
[edit]
user@switch# set applications application fcoe_app ether-type 0x8906
user@switch# set applications application iscsi_app protocol tcp destination-port 3260
```

13. Configure a DCBX application map to map the FCoE and iSCSI applications to the correct priorities:

```
[edit]
user@switch# set policy-options application-maps dcbx_iscsi_fcoe_app_map application
fcoe_app code-points [011 101]
user@switch# set policy-options application-maps dcbx_iscsi_fcoe_app_map application
iscsi_app code-points 100
```

14. Apply the application map to the interfaces so that DCBX exchanges FCoE application TLVs on the correct code points:

```
[edit]
user@switch# set protocols dcbx interface xe-0/0/31 application-map dcbx_iscsi_fcoe_app_map
user@switch# set protocols dcbx interface xe-0/0/32 application-map dcbx_iscsi_fcoe_app_map
user@switch# set protocols dcbx interface xe-0/0/33 application-map dcbx_iscsi_fcoe_app_map
user@switch# set protocols dcbx interface xe-0/0/34 application-map dcbx_iscsi_fcoe_app_map
```

Verification

IN THIS SECTION

- [Verifying the Forwarding Class Configuration | 610](#)
- [Verifying the Behavior Aggregate Classifier Configuration | 611](#)
- [Verifying the PFC Flow Control Configuration \(CNP\) | 613](#)
- [Verifying the Interface Configuration | 616](#)
- [Verifying the DCBX Application Configuration | 618](#)
- [Verifying the DCBX Application Map Configuration | 618](#)
- [Verifying the DCBX Application Protocol Exchange Interface Configuration | 619](#)

To verify the configuration and proper operation of the lossless forwarding classes and IEEE 802.1p priorities, perform these tasks:

Verifying the Forwarding Class Configuration

Purpose

Verify that the lossless forwarding classes `iscsi` and `fcoe1` have been created and that the default lossless forwarding class `fcoe` is still enabled for lossless transport.

Action

Show the forwarding class configuration by using the operational command `show class-of-service forwarding class`:

```
user@switch> show class-of-service forwarding-class
```

Forwarding class	ID	Queue	Policing priority	No-Loss
best-effort	0	0	normal	Disabled
fcoe	1	3	normal	Enabled
iscsi	2	4	normal	Enabled
network-control	3	7	normal	Disabled
fcoe1	4	5	normal	Enabled
mcast	8	8	normal	Disabled

Meaning

The `show class-of-service forwarding-class` command shows all of the forwarding classes. The command output shows that the `iscsi` and `fcoe1` forwarding classes are configured on output queues 4 and 5, respectively, with the no-loss packet drop attribute enabled.

Because we did not explicitly configure the default `fcoe` forwarding class, it remains in its default state (lossless configuration).

Verifying the Behavior Aggregate Classifier Configuration

Purpose

Verify that the four classifiers map the forwarding classes to the correct IEEE 802.1p code points (priorities) and packet loss priorities.

Action

List the classifiers configured to support lossless FCoE transport using the operational mode command `show class-of-service classifier`:

```
user@switch> show class-of-service classifier
```

Classifier: fcoe_p3_iscsi, Code point type: ieee-802.1, Index: 13915

Code point	Forwarding class	Loss priority
011	fcoe	low
100	iscsi	low

```

Classifier: fcoe_p5_iscsi, Code point type: ieee-802.1, Index: 62035
  Code point      Forwarding class      Loss priority
  100             iscsi                  low
  101             fcoe1                  low

Classifier: fcoe_p3_p5, Code point type: ieee-802.1, Index: 17774
  Code point      Forwarding class      Loss priority
  011             fcoe                  low
  101             fcoe1                  low

Classifier: iscsi_classifier, Code point type: ieee-802.1, Index: 31635
  Code point      Forwarding class      Loss priority
  100             iscsi                  low

```

Meaning

The `show class-of-service classifier` command shows the IEEE 802.1p code points and the loss priorities that are mapped to the forwarding classes in each classifier. The command output shows that there are four classifiers, `fcoe_p3_iscsi`, `fcoe_p5_iscsi`, `fcoe_p3_p5`, and `iscsi_classifier`.

Classifier `fcoe_p3_iscsi` maps code point 011 (priority 3) to default lossless forwarding class `fcoe` and a packet loss priority of `low`, and code point 100 (priority 4) to explicitly configured lossless forwarding class `iscsi`, and all other priorities to the best-effort forwarding class with a packet loss priority of `high`.

Classifier `fcoe_p5_iscsi` maps code point 100 to explicitly configured forwarding class `iscsi` and a packet loss priority of `low`, and code point 101 (priority 5) to explicitly configured lossless forwarding class `fcoe1` and a packet loss priority of `low`, and all other priorities to the best-effort forwarding class with a packet loss priority of `high`.

Classifier `fcoe_p3_p5` maps code point 011 to default lossless forwarding class `fcoe` and a packet loss priority of `low`, and maps code point 101 to explicitly configured lossless forwarding class `fcoe1` and a packet loss priority of `low`. The classifier maps all other priorities to the best-effort forwarding class with a packet loss priority of `high`.

Classifier `iscsi_classifier` maps code point 100 to explicitly configured forwarding class `iscsi` and a packet loss priority of `low`, and all other priorities to the best-effort forwarding class with a packet loss priority of `high`.

Verifying the PFC Flow Control Configuration (CNP)

Purpose

Verify that PFC is enabled on the correct input priorities and that flow control is configured on the correct output queues and priorities in each CNP.

Action

List the congestion notification profiles using the operational mode command `show class-of-service congestion-notification`:

```
user@switch> show class-of-service congestion-notification
Name: fcoe_p3_cnp, Index: 12037
Type: Input
Cable Length: 100 m
  Priority  PFC      MRU
  000      Disabled
  001      Disabled
  010      Disabled
  011      Enabled   2240
  100      Enabled   9216
  101      Disabled
  110      Disabled
  111      Disabled
Type: Output
  Priority  Flow-Control-Queues
  000
      0
  001
      1
  010
      2
  011
      3
  100
      4
  101
      5
  110
      6
  111
```


7

Name: fcoe_p3_p5_cnp, Index: 46484

Type: Input

Cable Length: 100 m

Priority	PFC	MRU
000	Disabled	
001	Disabled	
010	Disabled	
011	Enabled	2240
100	Disabled	
101	Enabled	2240
110	Disabled	
111	Disabled	

Type: Output

Priority	Flow-Control-Queues
011	
	3
101	
	5

Name: fcoe_p5_cnp, Index: 12133

Type: Input

Cable Length: 150 m

Priority	PFC	MRU
000	Disabled	
001	Disabled	
010	Disabled	
011	Disabled	
100	Enabled	9216
101	Enabled	2240
110	Disabled	
111	Disabled	

Type: Output

100	
	4
101	
	5

Name: iscsi_cnp, Index: 19342

Type: Input

Cable Length: 100 m

Priority	PFC	MRU
----------	-----	-----

```

000      Disabled
001      Disabled
010      Disabled
011      Disabled
100      Enabled      9216
101      Disabled
110      Disabled
111      Disabled
Type: Output
Priority  Flow-Control-Queues
000
      0
001
      1
010
      2
011
      3
100
      4
101
      5
110
      6
111
      7

```

Meaning

The `show class-of-service congestion-notification` command shows the input and output stanzas of the four CNPs.

For CNP `fcoe_p3_cnp`, the input stanza shows that PFC is enabled on IEEE 802.1p code point 011 (priority 3) with an MRU of 2240 bytes, and cable length of 100 meters. The input stanza also shows that PFC is enabled on code point 100 (priority 4) with the default MRU value of 9216 bytes. The CNP output stanza shows the default mapping of priorities to output queues because no explicit output CNP is configured.

NOTE: By default, only queues 3 and 4 are enabled respond to pause messages from the connected peer. For queue 3 to respond to pause messages, priority 3 (code point 011) must be

enabled for PFC in the input stanza. For queue 4 to respond to pause messages, priority 4 (code point 100) must be enabled for PFC in the input stanza. In this example, only queues 3 and 4 respond to pause messages from the connected peer on interfaces that use CNP fcoe_p3_cnp because the input stanza enables PFC only on priorities 3 and 4.

For CNP fcoe_p3_p5_cnp, the input stanza shows that PFC is enabled on code points 011 and 101 (priority 5), the MRU is 2240 bytes on both priorities, and the cable length is 100 meters. The CNP output stanza shows that output flow control is configured on queues 3 and 5 for code points 011 and 101, respectively.

For CNP fcoe_p5_cnp, the input stanza shows that PFC is enabled on code points 100 and 101. The MRU for code point 101 (FCoE traffic) is 2240 bytes and the MRU for code point 100 is 9216. The interface cable length is 150 meters. The CNP output stanza shows that output flow control is configured on queue 4 for code point 100 and on queue 5 for code point 101.

For CNP iscsi_cnp, the input stanza shows that PFC is enabled on code point 100, the MRU value is 9216 bytes, and the interface cable length is 100 meters. The CNP output stanza shows the default mapping of priorities to output queues because no explicit output CNP is configured.

Verifying the Interface Configuration

Purpose

Verify that the correct classifiers and congestion notification profiles are configured on the correct interfaces.

Action

List the ingress interfaces using the operational mode commands `show configuration class-of-service interfaces xe-0/0/31`, `show configuration class-of-service interfaces xe-0/0/32`, `show configuration class-of-service interfaces xe-0/0/33`, and `show configuration class-of-service interfaces xe-0/0/34`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/31
congestion-notification-profile fcoe_p3_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_p3_iscsi;
```

```
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/32
congestion-notification-profile fcoe_p5_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_p5_iscsi;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/33
congestion-notification-profile fcoe_p3_p5_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_p3_p5;
    }
}
```

```
user@switch> show configuration class-of-service interfaces xe-0/0/34
congestion-notification-profile iscsi_cnp;
unit 0 {
    classifiers {
        ieee-802.1 iscsi_classifier;
    }
}
```

Meaning

The `show configuration class-of-service interfaces xe-0/0/31` command shows that the congestion notification profile `fcoe_p3_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_p3_iscsi`.

The `show configuration class-of-service interfaces xe-0/0/32` command shows that the congestion notification profile `fcoe_p5_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_p5_iscsi`.

The `show configuration class-of-service interfaces xe-0/0/33` command shows that the congestion notification profile `fcoe_p3_p5_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_p3_p5`.

The `show configuration class-of-service interfaces xe-0/0/34` command shows that the congestion notification profile `iscsi_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `iscsi_classifier`.

Verifying the DCBX Application Configuration

Purpose

Verify that the DCBX applications for FCoE and iSCSI are configured.

Action

List the DCBX applications by using the configuration mode command `show applications`:

```
user@switch# show applications
application iscsi_app {
    protocol tcp;
    destination-port 3260;
}
application fcoe_app {
    ether-type 0x8906;
```

Meaning

The `show applications` configuration mode command shows all of the configured applications. The output shows that the application `iscsi_app` is configured with a protocol value of `tcp` and a destination port value of `3260`, and that the application `fcoe_app` is configured with an EtherType of `0x8906` (the correct EtherType for FCoE traffic).

Verifying the DCBX Application Map Configuration

Purpose

Verify that the application map is configured.

Action

List the application maps by using the configuration mode command `show policy-options application-maps`:

```
user@switch# show policy-options application-maps
dcbx-iscsi-fcoe-app-map {
    application iscsi_app code-points 100;
    application fcoe_app code-points [011 101];
}
```

Meaning

The `show policy-options application-maps` configuration mode command lists all of the configured application maps and the applications that belong to each application map. The output shows that there is one application map named `dcbx-iscsi-fcoe_app_map`. It consists of the application `iscsi_app` mapped to code point 100 and the application `fcoe_app` mapped to code points 011 and 101.

Verifying the DCBX Application Protocol Exchange Interface Configuration

Purpose

Verify that the application maps are applied to the correct interfaces.

Action

List the application maps on each interface using the configuration mode command `show protocols dcbx`:

```
user@switch# show protocols dcbx
interface xe-0/0/31.0 {
    application-map dcbx-iscsi-fcoe-app-map;
}
interface xe-0/0/32.0 {
    application-map dcbx-iscsi-fcoe-app-map;
}
interface xe-0/0/33.0 {
    application-map dcbx-iscsi-fcoe-app-map;
}
interface xe-0/0/34.0 {
```

```
application-map dcbx-iscsi-fcoe-app-map;
}
```

Meaning

The `show protocols dcbx configuration mode` command lists the application map association with interfaces. The output shows that all four interfaces use the application map `dcbx-iscsi-fcoe-app-map`.

RELATED DOCUMENTATION

Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface

Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch)

Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces

Example: Configuring DCBX Application Protocol TLV Exchange

Configuring CoS PFC (Congestion Notification Profiles)

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway

IN THIS SECTION

- [Priority Remapping Configuration | 621](#)
- [Configuration Rules | 622](#)
- [Fate Sharing | 623](#)

When the QFX Series acts as an FCoE-FC gateway, it connects an Ethernet network that carries Fibre Channel over Ethernet (FCoE) traffic to a Fibre Channel (FC) network. Ethernet interfaces connect to the FCoE network. Native FC interfaces (NP_Ports) connect to the FC network.

FCoE traffic typically uses IEEE 802.1p priority 3 (code point 011). The QFX Series default configuration maps priority 3 traffic to the FCoE forwarding class. If your FCoE network uses priority 3 for FCoE traffic, you do not need to remap priorities, because the default configuration maps priority 3 to the FCoE forwarding class. (But you do need to enable PFC on IEEE 802.1p code point 3 on the Ethernet interfaces to achieve lossless behavior.)

However, if the FCoE network uses a different IEEE 802.1p priority than priority 3 for FCoE traffic, then you can use priority remapping to classify FCoE traffic into a lossless forwarding class mapped to that priority (and classified to that priority on the FCoE Ethernet interfaces in the ingress classifier). You specify the lossless forwarding class used for the FCoE traffic by configuring a fixed classifier and applying it to the native FC (NP_Port) interface. All traffic received from the FC SAN on that NP_Port interface is classified into the forwarding class specified in the fixed classifier.

When native FC interfaces on the FCoE-FC gateway encapsulate incoming FC traffic in Ethernet to create FCoE frames, by default they assign IEEE 802.1p code point 011 to the FCoE traffic, forward the traffic internally to the gateway Ethernet interfaces, and then forward the traffic to the FCoE network. Setting a rewrite value for the IEEE 802.1p code point configures the gateway native FC interface to assign the rewrite value priority to the FCoE frames when the native FC interface forwards the FCoE frames to the gateway Ethernet interface. Instead of a priority of 3, the FCoE frames use the priority specified in the rewrite value.

You can configure one rewrite value for each local FCoE-FC gateway fabric. All of the native FC interfaces in a particular fabric must use the same rewrite value. Native FC interfaces that belong to different FCoE-FC gateway fabrics can use different rewrite values.

Priority Remapping Configuration

Native FC interfaces on an FCoE-FC gateway receive native FC traffic from the FC SAN and encapsulate it in Ethernet to create FCoE frames. Priority remapping enables you to map the encapsulated FC traffic (the FCoE traffic) to any IEEE 802.1p priority. (This is similar to the *rewrite rules* you can configure to remap forwarding classes to code points on Ethernet egress interfaces, but the rewrite takes place at the ingress FC interface so that the QFX Series uses the correct priority for FCoE traffic on the converged Ethernet network.)

To support lossless traffic flows, you must configure the remapped priority correctly on the native FC interfaces and also on the Ethernet interfaces that connect to the FCoE network. Achieving lossless behavior for FCoE traffic when you remap the FCoE priority requires configuring:

- A lossless forwarding class for FCoE traffic (or using the default *fcfe* forwarding class)
- A behavior aggregate (BA) classifier on the FCoE Ethernet interfaces to map the FCoE forwarding class to the IEEE 802.1p code points (priority) used for FCoE traffic on the FCoE network (the ingress classifier priority for the forwarding class must be the same as the rewrite value priority)

- A fixed classifier on the FCoE-FC gateway FC interface that maps all traffic from the FC network into the lossless FCoE forwarding class (the forwarding class must be lossless)
- A priority rewrite value that remaps the IEEE 802.1p code point on the FCoE-FC gateway FC interface to the priority used for FCoE traffic on the FCoE network
- An input congestion notification profile (CNP) to enable *priority-based flow control* (PFC) on the FCoE code point (the code point used as the rewrite value) at the Ethernet ingress interfaces

The ingress and egress configurations must match to achieve lossless behavior. The priority and the forwarding class specified in the BA classifier and in the CNP on the Ethernet ingress interfaces must match the fixed classifier and rewrite value on the FC interfaces. You must specify the same lossless FCoE forwarding class in each configuration and use the same IEEE 802.1p code point (priority) so that the FCoE traffic is properly classified into flows and so that those flows receive lossless treatment.

For example, if you configure a lossless forwarding class named *my_fcoe_fc* and your Ethernet network uses IEEE 802.1p priority 5 (code point 101) for FCoE traffic, then:

- The forwarding class configuration, the BA classifier, and the fixed classifier all specify *my_fcoe_fc* as the forwarding class
- The BA classifier, the input CNP, and the rewrite value all specify the IEEE 802.1p code point 101

Configuration Rules

The following configuration rules apply when you remap priorities on an FCoE-FC gateway:

- Each native FC interface (NP_Port) supports one IEEE 802.1p priority value. The interface rewrites the IEEE 802.1p code point of all incoming traffic on the interface to the rewrite value. (The FC interface uses either the default value of 3 or the rewrite value for all incoming traffic.)
- Ports in the same FCoE-FC gateway local fc-fabric must use the same rewrite value. For example, if ports fc-0/0/0 and fc-0/0/1 are in the same local FCoE-FC gateway fabric, they must use the same rewrite value. If you attempt to commit a configuration that uses different IEEE 802.1p priority rewrite values, the system returns a commit error.
- Ports in different FCoE-FC gateway local fc-fabrics can use different rewrite values. An example scenario is:
 - Interfaces fc-0/0/0 and fc-0/0/1 are in FCoE-FC gateway fc-fabric *my_fc_fab1*.
 - Interfaces fc-0/0/4 and fc-0/0/5 are in FCoE-FC gateway fc-fabric *my_fc_fab2*.

In this scenario, interfaces fc-0/0/0 and fc-0/0/1 must use the same rewrite value because they belong to the same local FC fabric on the gateway. Interfaces fc-0/0/4 and fc-0/0/5 also must use the same rewrite value because they belong to the same local FC fabric. However, the rewrite value

you use for interfaces fc-0/0/0 and fc-0/0/1 can be different than the rewrite value you use for interfaces fc-0/0/4 and fc-0/0/5 because the interfaces belong to different local FC fabrics.

- You can apply the rewrite value only to native FC interfaces; you cannot apply the rewrite value configuration to Ethernet interfaces.
- The forwarding class specified in the fixed classifier on the native FC interface must be a lossless forwarding class. You cannot apply a fixed classifier to a native FC interface unless the associated forwarding class is lossless. (The forwarding class must be one of the two default lossless forwarding classes, or you must explicitly configure the forwarding class with the *no-loss* drop attribute.)
- The lossless forwarding class and IEEE 802.1p priority configuration must match on the FCoE-FC gateway native FC interfaces and Ethernet interfaces:
 - The same IEEE 802.1p priority (code point) must be enabled for PFC on the Ethernet ingress interfaces, classified to the lossless forwarding class used in the native FC interface fixed classifier, and set as the rewrite value on the native FC interfaces.
 - The same lossless forwarding class must be used in the fixed classifier on the native FC interfaces and in the classifier configuration on the Ethernet interfaces.

Fate Sharing

To ensure that congestion on one interface does not affect the fate of traffic on a native FC interface on which you remap priorities, avoid fate sharing (different traffic flows receiving the same CoS treatment) configurations.

You can avoid fate sharing by ensuring that the remapping priority (code point) on the native FC interface is classified only to the forwarding class used in the fixed classifier on all other interfaces. For example, if you configure a fixed classifier on an FC interface that classifies all of the traffic into lossless forwarding class myfcoe1 and remaps the priority to priority 5 (IEEE 802.1p code point 101), then in all other classifier configurations on all other interfaces, priority 5 should always be classified to forwarding class myfcoe1. If you classify priority 6 on another interface to forwarding class myfcoe1, then congestion on priority 6 traffic affects priority 5 traffic unfairly.

RELATED DOCUMENTATION

[Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 528](#)

[Example: Configuring IEEE 802.1p Priority Remapping on an FCoE-FC Gateway | 624](#)

[Configuring CoS Fixed Classifier Rewrite Values for Native FC Interfaces \(NP_Ports\) | 156](#)

Example: Configuring IEEE 802.1p Priority Remapping on an FCoE-FC Gateway

IN THIS SECTION

- [Requirements | 624](#)
- [Overview | 625](#)
- [Configuration | 629](#)
- [Verification | 632](#)

FCoE traffic typically uses IEEE 802.1p priority 3 (code point 011). However, if your FCoE network uses a different IEEE 802.1p priority than priority 3 for FCoE traffic, then you can use priority remapping to classify FCoE traffic into a lossless forwarding class mapped to that priority. You specify the lossless forwarding class used for the FCoE traffic by configuring a fixed classifier and applying it to the native FC (NP_Port) interface. All traffic received from the FC SAN on that NP_Port interface is classified into the forwarding class specified in the fixed classifier.

When native FC interfaces on the FCoE-FC gateway encapsulate incoming FC traffic in Ethernet to create FCoE frames, by default they assign IEEE 802.1p code point 011 to the FCoE traffic, forward the traffic internally to the gateway Ethernet interfaces, and then forward the traffic to the FCoE network. Setting a rewrite value for the IEEE 802.1p code point configures the gateway native FC interface to assign the rewrite value priority to the FCoE frames when the native FC interface forwards the FCoE frames to the gateway Ethernet interface. Instead of a priority of 3, the FCoE frames use the priority specified in the rewrite value.

You can configure one rewrite value for each local FCoE-FC gateway fabric. All of the native FC interfaces in a particular fabric must use the same rewrite value. Native FC interfaces that belong to different FCoE-FC gateway fabrics can use different rewrite values.

This example shows how to configure FCoE priority remapping for a converged Ethernet network that uses priority 5 (IEEE code point 101) for FCoE traffic. If your network uses priority 3 for FCoE traffic, then you do not need to remap the FCoE priority, because the default configuration supports lossless FCoE transport on priority 3.

Requirements

This example uses the following hardware and software components:

- One Juniper Networks QFX3500 Switch

- Junos OS Release 12.3 or later for the QFX Series

Overview

IN THIS SECTION

- [Topology | 626](#)

Native FC interfaces on an FCoE-FC gateway receive native FC traffic from the FC SAN and encapsulate it in Ethernet to create FCoE frames. Priority remapping enables you to map the encapsulated FC traffic (the FCoE traffic) to any IEEE 802.1p priority.

To support lossless FCoE traffic flows, you must configure the remapped priority correctly on the native FC interfaces and also on the Ethernet interfaces that connect to the FCoE network. Achieving lossless behavior for FCoE traffic when you remap the FCoE priority requires configuring:

- A lossless forwarding class for FCoE traffic (or using the default `fcoe` forwarding class)
- A behavior aggregate (BA) classifier on the FCoE Ethernet interfaces to map the FCoE forwarding class to the IEEE 802.1p code points (priority) used for FCoE traffic on the FCoE network (the ingress classifier priority for the forwarding class must be the same as the rewrite value priority)
- A fixed classifier on the FCoE-FC gateway FC interface that maps all traffic from the FC network into the lossless FCoE forwarding class (the forwarding class must be lossless)
- A priority rewrite value that remaps the IEEE 802.1p code point on the FCoE-FC gateway FC interface to the priority used for FCoE traffic on the FCoE network
- An input congestion notification profile (CNP) to enable priority-based flow control (PFC) on the FCoE code point (the code point used as the rewrite value) at the Ethernet interface ingress and an output CNP to configure flow control to pause the correct output queue at the Ethernet interface egress

NOTE: Configuring or changing PFC on an interface blocks the entire port until the PFC change is completed. After a PFC change is completed, the port is unblocked and traffic resumes. Blocking the port stops ingress and egress traffic, and causes packet loss on all queues on the port until the port is unblocked.

- A DCBX application and application map on the Ethernet interface to support DCBX application TLV exchange for the lossless FCoE traffic on the FCoE priority

The priority specified in the BA classifier, CNP, and DCBX application map on the Ethernet ingress interfaces must match the priority specified in the fixed classifier and rewrite value configurations on the FC interfaces. You must specify the same lossless FCoE forwarding class in each configuration and use the same IEEE 802.1p code point (priority) so that the FCoE traffic is properly classified into flows and so that those flows receive lossless treatment.

Topology

This example shows how to configure priority remapping of FCoE traffic on one native FC interface (fc-0/0/2) connected to the FC SAN and on one Ethernet interface (xe-0/0/27) connected to the converged Ethernet (FCoE) network. Both the native FC interface and the Ethernet interface belong to the same local FC fabric on the FCoE-FC gateway.

The converged Ethernet network uses priority 5 (IEEE 802.1p code point 101) for FCoE traffic. The native FC interface on the FCoE-FC gateway receives FC traffic from the FC SAN. The native FC interface encapsulates the FC traffic in Ethernet to create FCoE frames, tags the frames with the IEEE 802.1p priority value 101, and then forwards the FCoE frames to the FCoE-FC gateway Ethernet interface. Because traffic marked with IEEE 802.1p priority 5 is mapped to a lossless FCoE forwarding class, the traffic receives lossless treatment. The Ethernet interface forwards the FCoE traffic on to the Ethernet network.

FCoE traffic (tagged with priority 5) arriving at the FCoE-FC gateway from the Ethernet network receives lossless treatment and is forwarded to the native FC interface. The native FC interface removes the Ethernet encapsulation from the FCoE frames and forwards the resulting native FC traffic to the FC SAN.

Figure 25 on page 626 shows the topology for this example, and Table 100 on page 627 shows the configuration components for this example.

Figure 25: Topology of the IEEE 802.1p Priority Remapping Example

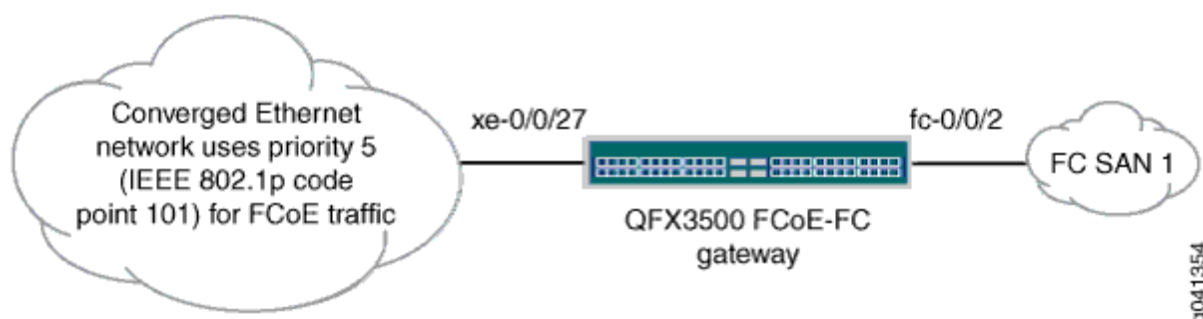


Table 100: Components of the IEEE 802.1p Priority Remapping Configuration Topology

Component	Settings
Hardware	QFX3500 switch
Forwarding class configuration	Name—fcoe1 Queue mapping—queue 5 Packet drop attribute—no-loss NOTE: The lossless forwarding class can be mapped to any output queue. However, because FCoE uses priority 5 in this example, matching that traffic to a forwarding class that uses queue 5 creates a configuration that is logical and easy to map because the priority and the queue are identified by the same number.
BA classifier (Ethernet interface)	Name—fcoe_gw_classifier Maps code point 101 (IEEE 802.1p priority 5) to the fcoe1 forwarding class and assigns traffic a packet loss priority of low. The classifier is applied to Ethernet interface xe-0/0/27.
Fixed classifier (native FC interface)	Forwarding class—fcoe1 The classifier is applied to native FC interface fc-0/0/2
Rewrite value	IEEE 802.1p code point—101 The rewrite value is applied to native FC interface fc-0/0/2
PFC configuration (CNP on Ethernet interface)	Name—fcoe1_p5_rewrite_cnp Input CNP code point—101 Output CNP code point—101 Output CNP flow control queue—5 Interface—xe-0/0/27

Table 100: Components of the IEEE 802.1p Priority Remapping Configuration Topology (Continued)

Component	Settings
DCBX application mapping	<p>Application name—myfcoe5</p> <p>Application ether-type—0x8906</p> <p>Application map name—myfcoe5_map</p> <p>Application map code points—101</p> <p>Interface—xe-0/0/27</p> <p>NOTE: LLDP and DCBX must be enabled on the interface. By default, LLDP and DCBX are enabled on all Ethernet interfaces.</p>

The priority used to identify FCoE traffic (5, IEEE 802.1p code point 101) is configured for lossless transport across the QFX device on interfaces xe-0/0/27 and fc-0/0/2, which belong to the same local FC fabric on the FCoE-FC gateway.

On the Ethernet interface, the classifier maps priority 5 to a lossless forwarding class (fcoe1), the input CNP enables PFC on incoming priority 5 traffic, and the output CNP enables output queue 5 to respond to pause messages received from the peer on traffic tagged with priority 5. On the native FC interface, FC traffic is remapped from priority 3 (the default mapping) to priority 5 and assigned to the same lossless forwarding class, fcoe1, because of the fixed classifier configuration. In this way, traffic tagged with priority 5 on interfaces xe-0/0/27 and fc-0/0/2 receives lossless treatment.

NOTE: To avoid fate sharing, ensure that the remapped priority is classified only to the forwarding class used in the fixed classifier on all other interfaces. For example, if you configure a fixed classifier on an FC interface that classifies all of the traffic into lossless forwarding class fcoe1 and remaps the priority to priority 5 (IEEE 802.1p code point 101), then in all other classifier configurations on all other interfaces, priority 5 should always be classified to forwarding class fcoe1. If you classify priority 6 on another interface to forwarding class fcoe1, then congestion on priority 6 traffic affects priority 5 traffic unfairly.

NOTE: This example does not include scheduling (bandwidth allocation) configuration or the local FC fabric configuration. This examples focuses only on priority remapping.

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 629](#)
- [Procedure | 629](#)

CLI Quick Configuration

To quickly configure IEEE 802.1p priority remapping on an FCoE-FC gateway, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level.

```
set class-of-service forwarding-classes class fcoe1 queue-num 5 no-loss
set class-of-service classifiers ieee-802.1 fcoe_gw_classifier forwarding-class fcoe1 loss-
priority low code-points 101
set class-of-service interfaces xe-0/0/27 unit 0 classifiers ieee-802.1 fcoe_gw_classifier
set class-of-service interfaces fc-0/0/2 forwarding-class fcoe1
set class-of-service interfaces fc-0/0/2 rewrite-value input ieee-802.1p code-point 101
set class-of-service congestion-notification-profile fcoe1_p5_rewrite_cnp input ieee-802.1 code-
point 101 pfc
set class-of-service congestion-notification-profile fcoe1_p5_rewrite_cnp output ieee-802.1 code-
point 101 flow-control-queue 5
set class-of-service interfaces xe-0/0/27 congestion-notification-profile fcoe1_p5_rewrite_cnp
set applications application myfcoe5 ether-type 0x8906
set policy-options application-maps myfcoe5_app_map application myfcoe5 code-points 101
set protocols dcbx interface xe-0/0/27 application-map myfcoe5_app_map
```

Procedure

Step-by-Step Procedure

To configure a lossless forwarding class for FCoE traffic, classify FCoE traffic into that forwarding class, configure a rewrite value on the native FC interface for the FCoE traffic, and enable PFC on the Ethernet interface, and configure DCBX application protocol TLV exchange for FCoE traffic:

1. Configure the lossless forwarding class (named `fcoe1` and mapped to output queue 5) for FCoE traffic that uses IEEE 802.1p priority 5:

```
[edit class-of-service]
user@switch# set forwarding-classes class fcoe1 queue-num 5 no-loss
```

2. Configure an ingress classifier named `fcoe_gw_classifier` to map the FCoE priority (IEEE 802.1p code point 101) to the lossless FCoE forwarding class (`fcoe1`):

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe_gw_classifier forwarding-class fcoe1 loss-priority low
code-points 101
```

3. Apply the classifier named `fcoe_gw_classifier` to Ethernet interface `xe-0/0/27`:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/27 unit 0 classifiers ieee-802.1 fcoe_gw_classifier
```

4. Configure the fixed classifier on the native FC interface, using the lossless FCoE forwarding class `fcoe1` (all traffic from the FC SAN is classified into the specified forwarding class). The traffic classified into this forwarding class is tagged with the priority value configured in the next step.

```
[edit class-of-service]
user@switch# set interfaces fc-0/0/2 forwarding-class fcoe1
```

5. Configure the rewrite value (IEEE 802.1p code point 101) applied to all incoming traffic from the FC SAN on the native FC interface. The rewrite value is the IEEE 802.1p priority that the encapsulated FCoE traffic classified into the `fcoe1` forwarding class uses on the converged Ethernet network.

```
[edit class-of-service]
user@switch# set interfaces fc-0/0/2 rewrite-value input ieee-802.1p code-point 101
```

6. Configure the input stanza of the CNP (named `fcoe1_p5_rewrite_cnp`) to enable PFC on the FCoE priority on the Ethernet interface:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe1_p5_rewrite_cnp input ieee-802.1 code-
point 101 pfc
```

7. Configure the output stanza of the CNP to enable output queue 5 to respond to pause messages received from the peer on traffic tagged with priority 5:

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe1_p5_rewrite_cnp output ieee-802.1
code-point 101 flow-control-queue 5
```

8. Apply the CNP named `fcoe1_p5_rewrite_cnp` to Ethernet interface `xe-0/0/27`:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/27 congestion-notification-profile fcoe1_p5_rewrite_cnp
```

9. Configure a DCBX application for FCoE to map to the Ethernet interface, so that DCBX can exchange application protocol TLVs on the correct (remapped) IEEE 802.1p FCoE priority:

```
[edit]
user@switch# set applications application myfcoe5 ether-type 0x8906
```

10. Configure a DCBX application map to map the FCoE application to the correct (remapped) IEEE 802.1p FCoE priority:

```
[edit]
user@switch# set policy-options application-maps myfcoe5_app_map application myfcoe5 code-
points 101
```

- 11. Apply the application map to the Ethernet interface so that DCBX exchanges FCoE application TLVs on the correct code point:

```
[edit]
user@switch# set protocols dcbx interface xe-0/0/27 application-map myfcoe5_app_map
```

Verification

IN THIS SECTION

- [Verifying the Forwarding Class Configuration | 632](#)
- [Verifying the Behavior Aggregate Classifier Configuration | 633](#)
- [Verifying the FC Interface Configuration \(Fixed Classifier, Rewrite Value\) | 634](#)
- [Verifying the Ethernet Interface PFC Configuration \(CNP\) | 634](#)
- [Verifying the Ethernet Interface Configuration | 635](#)
- [Verifying the DCBX Application Configuration | 636](#)
- [Verifying the DCBX Application Map Configuration | 637](#)
- [Verifying the DCBX Application Protocol Exchange Interface Configuration | 637](#)

To verify the configuration and proper operation of IEEE 802.1p priority remapping on an FCoE-FC gateway, perform these tasks:

Verifying the Forwarding Class Configuration

Purpose

Verify that the lossless forwarding class fcoe1 has been created.

Action

Show the forwarding class configuration by using the operational command `show class-of-service forwarding class`:

```
user@switch# show class-of-service forwarding-class
Forwarding class          ID      Queue  Policing priority  No-Loss
```

best-effort	0	0	normal	Disabled
fcoe	1	3	normal	Enabled
no-loss	2	4	normal	Enabled
network-control	3	7	normal	Disabled
fcoe1	4	5	normal	Enabled
mcast	8	8	normal	Disabled

Meaning

The `show class-of-service forwarding-class` command shows all of the forwarding classes. The command output shows that the `fcoe1` forwarding class is configured on output queue 5 with the `no-loss` packet drop attribute enabled.

Because we did not explicitly configure the default forwarding classes, they remain in their default state, including the lossless configuration of the `fcoe` and `no-loss` default forwarding classes.

Verifying the Behavior Aggregate Classifier Configuration

Purpose

Verify that the classifier maps the forwarding classes to the correct IEEE 802.1p code points (priorities) and packet loss priorities.

Action

List the classifier configured for priority remapping using the operational mode command `show class-of-service classifier name fcoe_gw_classifier`:

```
user@switch> show class-of-service classifier name fcoe_gw_classifier
Classifier: fcoe_gw_classifier, Code point type: ieee-802.1, Index: 13100
  Code point      Forwarding class      Loss priority
  101             fcoe1                  low
```

Meaning

The `show class-of-service classifier name fcoe_gw_classifier` command shows the IEEE 802.1p code points and the loss priorities that are mapped to the forwarding classes in the classifier. The command output shows that the classifier maps forwarding class `fcoe1` to IEEE 802.1p code point 101 (priority 5) with a packet loss priority of `low`.

Verifying the FC Interface Configuration (Fixed Classifier, Rewrite Value)

Purpose

Verify that the native FC interface (NP_Port) classifies incoming traffic into forwarding class fcoe1 and that the interface rewrite value is priority 5 (IEEE code point 101).

Action

Display the FC interface configuration using the operational mode command `show configuration class-of-service interfaces fc-0/0/2`:

```
user@switch> show configuration class-of-service interfaces fc-0/0/2
rewrite-value {
  input {
    ieee-802.1 {
      code-point {
        101;
      }
    }
  }
}
forwarding-class fcoe1;
```

Meaning

The `show configuration class-of-service interfaces fc-0/0/2` command shows that the rewrite value for incoming (input) traffic is IEEE 802.1p code point 101 (priority 5), and that the interface uses forwarding class fcoe1 as the fixed classifier for all incoming traffic.

Verifying the Ethernet Interface PFC Configuration (CNP)

Purpose

Verify that PFC is enabled on the correct priority (IEEE 802.1p code point 101) for lossless transport and that flow control is enabled on the correct output queue (queue 5) on the Ethernet interface.

Action

List the congestion notification profile using the operational mode command `show class-of-service congestion-notification fcoe1_p5_rewrite_cnp`:

```
user@switch> show class-of-service congestion-notification fcoe1_p5_rewrite_cnp
Name: fcoe1_p5_rewrite_cnp, Index: 7061
Type: Input
Cable Length: 100 m
  Priority    PFC        MRU
  000        Disabled
  001        Disabled
  010        Disabled
  011        Disabled
  100        Disabled
  101        Enabled    2500
  110        Disabled
  111        Disabled
Type: Output
  Priority    Flow-Control-Queues
  101
      5
```

Meaning

The `show class-of-service congestion-notification fcoe1_p5_rewrite_cnp` command shows the input and output stanzas of the CNP. The input stanza shows that PFC is enabled on IEEE 802.1p code point 101 (priority 5). The input stanza also shows that the CNP uses the default values of 100 meters for the cable length value and 2500 bytes for the maximum receive unit (MRU) value.

The output stanza shows that flow control is enabled on output queue 5 for IEEE 802.1p priority code point 101 (priority 5).

Verifying the Ethernet Interface Configuration

Purpose

Verify that the classifier `fcoe_gw_classifier` and the congestion notification profile `fcoe1_p5_rewrite_cnp` are configured on Ethernet interface `xe-0/0/27`.

Action

List the ingress interfaces using the operational mode command `show configuration class-of-service interfaces xe-0/0/27`:

```
user@switch> show configuration class-of-service interfaces xe-0/0/27
congestion-notification-profile fcoe1_p5_rewrite_cnp;
unit 0 {
    classifiers {
        ieee-802.1 fcoe_gw_classifier;
    }
}
```

Meaning

The `show configuration class-of-service interfaces xe-0/0/27` command shows that the congestion notification profile `fcoe1_p5_rewrite_cnp` is configured on the interface, and that the IEEE 802.1p classifier associated with the interface is `fcoe_gw_classifier`.

Verifying the DCBX Application Configuration

Purpose

Verify that the DCBX application named `myfcoe5` for FCoE is configured.

Action

List the DCBX applications by using the configuration mode command `show applications`:

```
user@switch# show applications
application myfcoe5 {
    ether-type 0x8906;
}
```

Meaning

The `show applications` configuration mode command shows all of the configured applications. The output shows that the application `myfcoe5` is configured with an EtherType of `0x8906` (the correct EtherType for FCoE traffic).

Verifying the DCBX Application Map Configuration

Purpose

Verify that the application map `myfcoe5_app_map` is configured.

Action

List the application map by using the configuration mode command `show policy-options application-maps`:

```
user@switch# show policy-options application-maps
myfcoe5_app_map {
    application myfcoe5 code-points 101;
}
```

Meaning

The `show policy-options application-maps` configuration mode command lists all of the configured application maps and the applications that belong to each application map. The output shows that there is one application map, `myfcoe5_app_map`, which consists of the application named `myfcoe5` mapped to IEEE 802.1p code point 101 (priority 5).

Verifying the DCBX Application Protocol Exchange Interface Configuration

Purpose

Verify that the application map is applied to the correct interface (`xe-0/0/27`).

Action

List the application maps using the configuration mode command `show protocols dcbx`:

```
user@switch# show protocols dcbx
interface xe-0/0/27.0 {
    application-map myfcoe5_app_map;
}
```


Meaning

The `show protocols dcbx configuration mode` command lists the application map association with interfaces. The output shows that interface `xe-0/0/27` uses application map `myfcoe5_app_map`.

RELATED DOCUMENTATION

[Example: Configuring DCBX Application Protocol TLV Exchange | 661](#)

[Example: Configuring Unicast Classifiers | 86](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 507](#)

[Configuring CoS Fixed Classifier Rewrite Values for Native FC Interfaces \(NP_Ports\) | 156](#)

[Understanding CoS IEEE 802.1p Priority Remapping on an FCoE-FC Gateway | 620](#)

Understanding DCBX

IN THIS SECTION

- [DCBX Basics | 639](#)
- [DCBX Modes and Support | 640](#)
- [DCBX Attribute Types | 643](#)
- [DCBX Application Protocol TLV Exchange | 644](#)
- [DCBX and PFC | 645](#)
- [DCBX and ETS | 645](#)

Data Center Bridging Capability Exchange protocol (DCBX) is an extension of Link Layer Data Protocol (LLDP). If you disable LLDP on an interface, that interface cannot run DCBX. If you attempt to enable DCBX on an interface on which LLDP is disabled, the configuration commit operation fails. Data center bridging (DCB) devices use DCBX to exchange configuration information with directly connected peers.



Video: [What is DCBX Protocol?](#)

This topic describes:

DCBX Basics

DCBX can:

- Discover the DCB capabilities of peers.
- Detect DCB feature misconfiguration or mismatches between peers.
- Configure DCB features on peers.

You can configure DCBX operation for *priority-based flow control* (PFC), Layer 2 and Layer 4 applications such as FCoE and iSCSI, and ETS. DCBX is enabled or disabled on a per-interface basis.

NOTE: QFX5200 and QFX5210 switches do not support enhanced transmission selection (ETS) hierarchical scheduling. Use port scheduling to manage bandwidth on these switches.

By default, for PFC and ETS, DCBX automatically negotiates administrative state and configuration with each interface's connected peer. To enable DCBX negotiation for applications, you must configure the applications, map them to IEEE 802.1p code points in an application map, and apply the application map to interfaces.

The FCoE application only needs to be included in an application map when you want an interface to exchange type, length, and values (TLVs) for other applications in addition to FCoE. If FCoE is the only application you want an interface to advertise, then you do not need to use an application map. For ETS, DCBX pushes the switch configuration to peers if they are set to learn the configuration from the switch (unless you disable sending the ETS recommendation TLV on interfaces in IEEE DCBX mode).

You can override the default behavior for PFC, for ETS, or for all applications mapped to an interface by turning off autonegotiation to force an interface to enable or disable that feature. You can also disable DCBX autonegotiation for applications on an interface by excluding those applications from the application map you apply to that interface or by deleting the application map from the interface.

The default autonegotiation behavior for applications that are mapped to an interface is:

- DCBX is enabled on the interface if the connected peer device also supports DCBX.
- DCBX is disabled on the interface if the connected peer device does not support DCBX.

During negotiation of capabilities, the switch can push the PFC configuration to an attached peer if the peer is configured as "willing" to learn the PFC configuration from other peers. The Juniper Networks switch does not support self autoprovisioning and does not change its configuration during autonegotiation to match the peer configuration. (The Juniper switch is not "willing" to learn the PFC configuration from peers.)

NOTE: When a port with DCBX enabled begins to exchange type, length, and value (TLV) entries, optional LLDP TLVs on that port are not advertised to neighbors, so that the switch can interoperate with a wider variety of converged network adapters (CNAs) and Layer 2 switches that support DCBX.

DCBX Modes and Support

This section describes DCBX support:

DCBX Modes (Versions)

The two most common DCBX modes are supported:

- IEEE DCBX—The newest DCBX version. Different TLVs have different subtypes (for example, the subtype for the ETS configuration TLV is 9); the IEEE DCBX Organizationally Unique Identifier (OUI) is 0x0080c2.
- DCBX version 1.01—The Converged Enhanced Ethernet (CEE) version of DCBX. It has a subtype of 2 and an OUI of 0x001b21.

IEEE DCBX and DCBX version 1.01 differ mainly in frame format. DCBX version 1.01 uses one TLV that includes all DCBX attribute information, which is sent as sub-TLVs. IEEE DCBX uses a unique TLV for each DCB attribute.

NOTE: The switch does not support pre-CEE (pre-DCB) DCBX versions. Unsupported older versions of DCBX have a subtype of 1 and an OUI of 0x001b21. The switch drops LLDP frames that contain pre-CEE DCBX TLVs.

Table 101 on page 640 summarizes the differences between IEEE DCBX and DCBX version 1.01, including show command output:

Table 101: Summary of Differences Between IEEE DCBX and DCBX Version 1.01

Characteristic	IEEE DCBX	DCBX Version 1.01
OUI	0x0080c2	0x001b21

Table 101: Summary of Differences Between IEEE DCBX and DCBX Version 1.01 (Continued)

Characteristic	IEEE DCBX	DCBX Version 1.01
Frame Format	Sends a separate, unique TLV for each DCBX attribute. For example, IEEE DCBX uses separate TLVs for ETS, PFC, and each application. Configuration and Recommendation information is sent in different TLVs	Sends one TLV that includes all DCBX attribute information organized in sub-TLVs. The “willing” bit determines whether or not an interface can change its configuration to match the connected peer.
Symmetric/asymmetric configuration with peer	Asymmetric or symmetric	Symmetric only
Differences in the <code>show dcbx interface interface-name</code> operational command	<ul style="list-style-type: none"> • Synchronization information is not shown because symmetric configuration is not required. • Operational state information is not shown because the operational states do not have to be symmetric. • TLV type is shown because unique TLVs are sent for each DCBX attribute. • ETS peer Configuration TLV and Recommendation TLV information is shown separately because they are different TLVs. 	<ul style="list-style-type: none"> • Synchronization information is shown because symmetric configuration is required. • Operational state information is shown because the operational states do have to be symmetric. • TLV type is not shown because one TLV is used for all attribute information. • Recommendation TLV is not sent (DCBX Version 1.01 uses the “willing” bit to determine whether or not an interface uses the peer interface configuration).

You can configure interfaces to use the following DCBX modes:

- IEEE DCBX—The interface uses IEEE DCBX regardless of the configuration on the connected peer.
- DCBX version 1.01—The interface uses DCBX version 1.01 regardless of the configuration on the connected peer.
- Autonegotiation—The interface automatically negotiates with the connected peer to determine the DCBX version the peers use. Autonegotiation is the default DCBX mode.

If you configure a DCBX mode on an interface, the interface ignores DCBX protocol data units (PDUs) it receives from the connected peer if the PDUs do not match the DCBX version configured on the

interface. For example, if you configure an interface to use IEEE DCBX and the connected peer sends DCBX version 1.01 LLDP PDUs, the interface ignores the version 1.01 PDUs. If you configure an interface to use DCBX version 1.01 and the peer sends IEEE DCBX LLDP PDUs, the interface ignores the IEEE DCBX PDUs.

NOTE: On interfaces that use the IEEE DCBX mode, the `show dcbx neighbors interface interface-name` operational command does not include application, PFC, or ETS operational state in the output.

Autonegotiation

Autonegotiation is the default DCBX mode. Each interface automatically negotiates with its connected peer to determine the DCBX version that both interfaces use to exchange DCBX information.

When an interface connects to its peer interface, the interface advertises IEEE DCBX TLVs to the peer. If the interface receives one IEEE DCBX PDU from the peer, the interface sets the DCBX mode as IEEE DCBX. If the interface receives three DCBX version 1.01 TLVs from the peer, the interface sets DCBX version 1.01 as the DCBX mode.

Autonegotiation works slightly differently on standalone switches compared to QFabric systems:

- Standalone switches—When an interface connects to its peer interface, the interface advertises IEEE DCBX TLVs to the peer. If the interface receives an IEEE DCBX TLV from the peer, the interface sets IEEE DCBX as the DCBX mode. If the interface receives three consecutive DCBX version 1.01 TLVs from the peer, the interface sets DCBX version 1.01 as the DCBX mode.
- QFabric system—When an interface connects to its peer interface, the interface advertises DCBX version 1.01 TLVs to the peer. If the interface receives an IEEE DCBX TLVs from the peer, the interface sets IEEE DCBX as the DCBX mode. If the interface receives three consecutive DCBX version 1.01 TLVs from the peer, the interface retains DCBX version 1.01 as the DCBX mode.

NOTE: If the link flaps or the LLDP process restarts, the interface starts the autonegotiation process again. The interface does not use the last received DCBX communication mode.

CNA Support for DCBX Modes

Different CNA vendors support different versions and capabilities of DCBX. The DCBX configuration you use on switch interfaces depends on the DCBX features that the CNAs in your network support.

Interface Support for DCBX

You can configure DCBX on 10-Gigabit Ethernet interfaces and on link aggregation group (LAG) interfaces whose member interfaces are all 10-Gigabit Ethernet interfaces.

DCBX Attribute Types

DCBX has three attribute types:

- **Informational**—These attributes are exchanged using LLDP, but do not affect DCBX state or operation; they only communicate information to the peer. For example, application priority TLVs are informational TLVs.
- **Asymmetric**—The values for these types of attributes do not have to be the same on the connected peer interfaces. Peers exchange asymmetric attributes when the attribute values can differ on each peer interface. The peer interface configurations might match or they might differ. For example, ETS Configuration and Recommendation TLVs are asymmetric TLVs.
- **Symmetric**—The intention is that the values for these types of attributes should be the same on both of the connected peer interfaces. Peer interfaces exchange symmetric attributes to ensure symmetric DCBX configuration for those attributes. For example, PFC Configuration TLVs are symmetric TLVs.

The following sections describe asymmetric and symmetric DCBX attributes:

Asymmetric Attributes

DCBX passes asymmetric attributes between connected peer interfaces to communicate parameter information about those attributes (features). The resulting configuration for an attribute might be different on each peer, so the parameters configured on one interface might not match the parameters on the connected peer interface.

There are two types of asymmetric attribute TLVs:

- **Configuration TLV**—Configuration TLVs communicate the current operational state and the state of the “willing” bit. The “willing” bit communicates whether or not the interface is willing to accept and use the configuration from the peer interface. If an interface is “willing,” the interface uses the configuration it receives from the peer interface. (The peer interface configuration can override the configuration on the “willing” interface.) If an interface is “not willing,” the configuration on the interface cannot be overridden by the peer interface configuration.
- **Recommendation TLV**—Recommendation TLVs communicate the parameters the interface recommends that the connected peer interface should use. When an interface sends a Recommendation TLV, if the connected peer is “willing,” the connected peer changes its configuration to match the parameters in the Recommendation TLV.

Symmetric Attributes

DCBX passes symmetric attributes between connected peer interfaces to communicate parameter information about those attributes (features), with the objective that both interfaces should use the same configuration. The intent is that the parameters configured on one interface should match the parameters on the connected peer interface.

There is one type of symmetric attribute TLV, the Configuration TLV. As with asymmetric attributes, symmetric attribute Configuration TLVs communicate the current operational state and the state of the “willing” bit. “Willing” interfaces use the peer interface parameter values for the attribute. (The attribute configuration of the peer overrides the configuration on the “willing” interface.)

DCBX Application Protocol TLV Exchange

DCBX advertises the switch’s capabilities for Layer 2 applications such as FCoE and Layer 4 applications such as iSCSI:

Application Protocol TLV Exchange

For all applications, DCBX advertises the application’s state and IEEE 802.1p code points on the interfaces to which the application is mapped. If an application is not mapped to an interface, that interface does not advertise the application’s TLVs. There is an exception for FCoE application protocol TLV exchange when FCoE is the only application you want DCBX to advertise on an interface.

FCoE Application Protocol TLV Exchange

Protocol TLV exchange for the FCoE application depends on whether FCoE is the only application you want the interface to advertise or whether you want the interface to exchange other application TLVs in addition to FCoE TLVs.

If FCoE is the only application you want DCBX to advertise on an interface, DCBX exchanges FCoE application protocol TLVs by default if the interface:

- Carries FCoE traffic (traffic mapped by CoS configuration to the FCoE forwarding class)
- Has a congestion notification profile with PFC enabled on the FCoE priority (IEEE 802.1p code point)
- Does *not* have an application map

NOTE: If no CoS configuration for FCoE is mapped to an interface, that interface does not exchange FCoE application protocol TLVs.

If you want DCBX to advertise FCoE and other applications on an interface, you must specify all of the applications, including FCoE, in an application map, and apply the application map to the desired interfaces.

NOTE: If an application map is applied to an interface, the FCoE application must be explicitly configured in the application map, or the interface does not exchange FCoE TLVs.

When DCBX advertises the FCoE application, it advertises the FCoE state and IEEE 802.1p code points. If a peer device connected to a switch interface does not support FCoE, DCBX uses autonegotiation to mark the interface as “FCoE down,” and FCoE is disabled on that interface.

Disabling Application Protocol TLV Exchange

To disable DCBX application protocol exchange for all applications on an interface, issue the `set protocols dcbx interface interface-name applications no-auto-negotiation` command.

You can also disable DCBX application protocol exchange for applications on an interface by deleting the application map from the interface, or by deleting a particular application from the application map. However, when you delete an application from an application map, the application protocol is no longer exchanged on any interface which uses that application map.

DCBX and PFC

After you enable PFC on a switch interface, DCBX uses autonegotiation to control the operational state of the PFC functionality.

If the peer device connected to the interface supports PFC and is provisioned compatibly with the switch, DCBX sets the PFC operational state to enabled. If the peer device connected to the interface does not support PFC or is not provisioned compatibly with the switch, DCBX sets the operational state to disabled. (PFC must be symmetrical.)

If the peer advertises that it is “willing” to learn its PFC configuration from the switch, DCBX pushes the switch’s PFC configuration to the peer and does not check the peer’s administrative state.

You can manually override DCBX control of the PFC operational state on a per-interface basis by disabling autonegotiation. If you disable autonegotiation on an interface on which you have configured PFC, then PFC is enabled on that interface regardless of the peer configuration. To disable PFC on an interface, do not configure PFC on that interface.

DCBX and ETS

This section describes:

Default DCBX ETS Advertisement

If you do not configure ETS on an interface, the switch automatically creates a default priority group that contains all of the priorities (forwarding classes, which represent output queues) and assigns 100 percent of the port output bandwidth to that priority group. The default priority group is transparent. It does not appear in the configuration and is used for DCBX advertisement. DCBX advertises the default priority group, its priorities, and the assigned bandwidth.

If you configure ETS on an interface, DCBX advertises:

- Each priority group on the interface
- The priorities in each priority group
- The bandwidth properties of each priority group and priority

Any priority on that interface that is not part of an explicitly configured priority group (forwarding class set) is assigned to the automatically generated default priority group and receives no bandwidth. If you configure ETS on an interface, every forwarding class (priority) on that interface for which you want to forward traffic must belong to a forwarding class set (priority group).

ETS Advertisement and Peer Configuration

DCBX does not control the switch's ETS (hierarchical scheduling) operational state. If the connected peer is configured as "willing," DCBX pushes the switch's ETS configuration to the switch's peers if the ETS Recommendation TLV is enabled (it is enabled by default). If the peer does not support ETS or is not consistently provisioned with the switch, DCBX does not change the ETS operational state on the switch. The ETS operational state remains enabled or disabled based only on the switch hierarchical scheduling configuration and is enabled by default.

When ETS is configured, DCBX advertises the priority groups, the priorities in the priority groups, and the bandwidth configuration for the priority groups and priorities. Any priority (essentially a forwarding class or queue) that is not part of a priority group has no scheduling properties and receives no bandwidth.

You can manually override whether DCBX advertises the ETS state to the peer on a per-interface basis by disabling autonegotiation. This does not affect the ETS state on the switch or on the peer, but it does prevent the switch from sending the Recommendation TLV or the Configuration TLV to the connected peer. To disable ETS on an interface, do not configure priority groups (forwarding class sets) on the interface.

ETS Recommendation TLV

The ETS Recommendation TLV communicates the ETS settings that the switch wants the connected peer interface to use. If the peer interface is "willing," it changes its configuration to match the

configuration in the ETS Recommendation TLV. By default, the switch interfaces send the ETS Recommendation TLV to the peer. The settings communicated are the egress ETS settings defined by configuring hierarchical scheduling on the interface.

We recommend that you use the same ETS settings on the connected peer that you use on the switch interface and that you leave the ETS Recommendation TLV enabled. However, on interfaces that use IEEE DCBX as the DCBX mode, if you want an asymmetric configuration between the switch interface and the connected peer, you can disable the ETS Recommendation TLV by including the `no-recommendation-tlv` statement at the `[edit protocols dcbx interface interface-name enhanced-transmission-selection]` hierarchy level.

NOTE: You can disable the ETS Recommendation TLV only when the DCBX mode on the interface is IEEE DCBX. Disabling the ETS Recommendation TLV has no effect if the DCBX mode on the interface is DCBX version 1.01. (IEEE DCBX uses separate application attribute TLVs, but DCBX version 1.01 sends all application attributes in the same TLV and uses sub-TLVs to separate the information.)

If you disable the ETS Recommendation TLV, the switch still sends the ETS Configuration TLV to the connected peer. The result is that the connected peer is informed about the switch DCBX ETS configuration, but even if the peer is “willing,” the peer does not change its configuration to match the switch configuration. This is asymmetric configuration—the two interfaces can have different parameter values for the ETS attribute.

For example, if you want a CNA connected to a switch interface to have different bandwidth allocations than the switch ETS configuration, you can disable the ETS Recommendation TLV and configure the CNA for the desired bandwidth. The switch interface and the CNA exchange configuration parameters, but the CNA does not change its configuration to match the switch interface configuration.

Release History Table

Release	Description
21.2R1EVO	PTX10008 routers support DCBX and PFC.

RELATED DOCUMENTATION

Understanding DCBX Application Protocol TLV Exchange
Understanding DCB Features and Requirements
Understanding CoS Flow Control (Ethernet PAUSE and PFC)
Understanding CoS Hierarchical Port Scheduling (ETS)

Understanding CoS Port Schedulers on QFX Switches

Understanding FCoE

Configuring the DCBX Mode

Configuring DCBX Autonegotiation

Disabling the ETS Recommendation TLV

Example: Configuring DCBX Application Protocol TLV Exchange

Configuring the DCBX Mode

You can configure the DCBX mode that an interface uses to communicate with the connected peer. Three DCBX modes are supported:

- **Autonegotiation**—The interface negotiates with the connected peer to determine the DCBX mode. This is the default DCBX mode.
- **IEEE DCBX**—The interface uses IEEE DCBX type, length, and value (TLV) to exchange DCBX information with the connected peer. QFX3500 Node devices come up with IEEE DCBX enabled by default and then autonegotiate with the connected peer to determine the final DCBX mode.
- **DCBX Version 1.01**—The interface uses Converged Enhanced Ethernet (CEE) DCBX version 1.01 TLVs to exchange DCBX information with the connected peer. QFabric system Node devices other than QFX3500 switches come up with DCBX version 1.01 enabled by default and then autonegotiate with the connected peer to determine the final DCBX mode.

NOTE: Pre-CEE (pre-DCB) versions of DCBX such as DCBX version 1.00 are not supported. If an interface receives an LLDP frame with pre-CEE DCBX TLVs, the system drops the frame.

Configure the DCBX mode by specifying the mode for one interface or for all interfaces.

- To configure the DCBX mode, specify the interface and the mode:

```
[edit protocols dcbx]
user@switch# set interface interface-name dcbx-version (auto-negotiate | ieee-dcbx | dcbx-
version-1.01)
```

For example, to configure DCBX version 1.01 on interface xe-0/0/21:

```
user@switch# set protocols dcbx interface xe-0/0/21 dcbx-version dcbx-version-1.01
```

To configure IEEE DCBX on all interfaces:

```
user@switch# set protocols dcbx interface all dcbx-version ieee-dcbx
```

RELATED DOCUMENTATION

Configuring DCBX Autonegotiation

Disabling the ETS Recommendation TLV

Understanding DCBX

Understanding DCBX Application Protocol TLV Exchange

show dcbx neighbors

Configuring DCBX Autonegotiation

Data Center Bridging Capability Exchange protocol (DCBX) discovers the data center bridging (DCB) capabilities of peers by exchanging feature configuration information. DCBX also detects feature misconfiguration and mismatches, and can configure DCB on peers. DCBX is an extension of the Link Layer Discovery Protocol (LLDP), and LLDP must remain enabled on every interface for which you want to use DCBX. If you attempt to enable DCBX on an interface on which LLDP is disabled, the configuration commit operation fails.

NOTE: LLDP and DCBX are enabled by default on all interfaces.

The switch supports DCBX autonegotiation for:

- Priority-based flow control (PFC) configuration
- Layer 2 and Layer 4 applications such as Fibre Channel over Ethernet (FCoE) and Internet Small Computer System Interface (iSCSI)
- Enhanced transmission selection (ETS) advertisement

DCBX autonegotiation is configured on a per-interface basis for each supported feature or application. The PFC and application DCBX exchanges use autonegotiation by default. The default autonegotiation behavior is:

- DCBX is enabled on the interface if the connected peer device also supports DCBX.
- DCBX is disabled on the interface if the connected peer device does not support DCBX.

You can override the default behavior for each feature by turning off autonegotiation to force an interface to enable or disable the feature.

Autonegotiation of ETS means that when ETS is enabled on an interface (priority groups are configured), the interface advertises its ETS configuration to the peer device. In this case, priorities (forwarding classes) that are not part of a priority group (forwarding class set) receive no bandwidth and are advertised in an automatically generated default forwarding class. If ETS is not enabled on an interface (no priority groups are configured), all of the priorities are advertised in one automatically generated default priority group that receives 100 percent of the port bandwidth.

Disabling ETS autonegotiation prevents the interface from sending the Recommendation TLV or the Configuration TLV to the connected peer.

On interfaces that use IEEE DCBX mode to exchange DCBX parameters, you can disable autonegotiation of the ETS Recommendation TLV to the peer if you want an asymmetric ETS configuration between the peers. DCBX still exchanges the ETS Configuration TLV if you disable the ETS Recommendation TLV.

Autonegotiation of PFC means that when PFC is enabled on an interface, if the peer device connected to the interface supports PFC and is provisioned compatibly with the switch, DCBX sets the PFC operational state to enabled. If the peer device connected to the interface does not support PFC or is not provisioned compatibly with the switch, DCBX sets the operational state to disabled.

In addition, if the peer advertises that it is “willing” to learn its PFC configuration from the switch, DCBX pushes the switch’s PFC configuration to the peer and does not check the peer’s administrative state. The switch does not learn PFC configuration from peers (the switch does not advertise its state as “willing”).

Disabling PFC autonegotiation prevents the interface from exchanging PFC configuration information with the peer. It forces the interface to enable PFC if PFC is configured on the interface or to disable PFC if PFC is not configured on the interface. If you disable PFC autonegotiation, the assumption is that the peer is also configured manually.

Autonegotiation of applications depends on whether or not you apply an application map to an interface. If you apply an application map to an interface, the interface autonegotiates DCBX for each application in the application map. PFC must be enabled on the FCoE priority (the FCoE IEEE 802.1p code point) for the interface to advertise the FCoE application. The interface only advertises applications that are included in the application map.

For example, if you apply an application map to an interface and the application map does not include the FCoE application, then that interface does not perform DCBX advertisement of FCoE.

If you do not apply an application map to an interface, DCBX does not advertise applications on that interface, with the exception of FCoE, which is handled differently than other applications.

NOTE: If you do not apply an application map to an interface, the interface performs autonegotiation of FCoE if the interface carries traffic in the FCoE forwarding class and also has PFC enabled on the FCoE priority. On such interfaces, if DCBX detects that the peer device connected to the interface supports FCoE, the switch advertises its FCoE capability and IEEE 802.1p code point on that interface. If DCBX detects that the peer device connected to the interface does not support FCoE, DCBX marks that interface as “FCoE down” and disables FCoE on the interface.

When DCBX marks an interface as “FCoE down,” the behavior of the switch depends on how you use it in the network:

- When the switch acts as an FCoE transit switch, the interface drops all of the FIP packets it receives. In addition, FIP packets received from an FCoE forwarder (FCF) are not forwarded to interfaces marked as “FCoE down.”
- When the switch acts as an FCoE-FC gateway (only switches that support native Fibre Channel interfaces), it does not send or receive FCoE Initialization Protocol (FIP) packets.

Disabling autonegotiation prevents the interface from exchanging application information with the peer. In this case, the assumption is that the peer is also configured manually.

To disable DCBX autonegotiation of PFC, applications (including FCoE), and ETS using the CLI:

1. Turn off autonegotiation for PFC.

```
[edit]
user@switch# set protocols dcbx interface interface-name priority-flow-control no-auto-
negotiation
```

2. Turn off autonegotiation for applications.

```
[edit]
user@switch# set protocols dcbx interface interface-name applications no-auto-negotiation
```

3. Turn off autonegotiation for ETS.

```
[edit]
user@switch# set protocols dcbx interface interface-name enhanced-transmission-selection no-
auto-negotiation
```

To disable autonegotiation of the ETS Recommendation TLV so that DCBX exchanges only the ETS Configuration TLV:

- ```
[edit protocols dcbx interface interface-name]
user@switch# set enhanced-transmission-selection no-recommendation-tlv
```

## RELATED DOCUMENTATION

*Example: Configuring DCBX Application Protocol TLV Exchange*

*Example: Configuring CoS PFC for FCoE Traffic*

*Disabling the ETS Recommendation TLV*

*Understanding DCBX Application Protocol TLV Exchange*

## Understanding DCBX Application Protocol TLV Exchange

### IN THIS SECTION

- [Applications | 653](#)
- [Application Maps | 654](#)
- [Classifying and Prioritizing Application Traffic | 655](#)
- [Enabling Interfaces to Exchange Application Protocol Information | 656](#)
- [Disabling DCBX Application Protocol Exchange | 656](#)

Data Center Bridging Capability Exchange protocol (DCBX) discovers the data center bridging (DCB) capabilities of connected peers. DCBX also advertises the capabilities of applications on interfaces by exchanging application protocol information through application type, length, and value (TLV) elements.

DCBX is an extension of Link Layer Discovery Protocol (LLDP). LLDP must remain enabled on every interface on which you want to use DCBX.

**NOTE:** LLDP and DCBX are enabled by default on all interfaces.

Setting up application protocol exchange consists of:

- Defining applications
- Mapping the applications to IEEE 802.1p code points in an *application map*
- Configuring classifiers to prioritize incoming traffic and map the incoming traffic to the application by the traffic code points
- Applying the application maps and classifiers to interfaces

You need to explicitly define the applications that you want an interface to advertise. The FCoE application is a special case (see ["Applications" on page 653](#)) and only needs to be defined on an interface if you want DCBX to exchange application protocol TLVs for other applications in addition to FCoE on that interface.

You also need to explicitly map all of the defined applications that you want an interface to advertise to IEEE 802.1p code points in an application map. The FCoE application is a special case that only requires inclusion in an application map when you want an interface to use DCBX for other applications in addition to FCoE, as described later in this topic (see ["Application Maps" on page 654](#)).

This topic describes:

## Applications

Before an interface can exchange application protocol information, you need to define the applications that you want to advertise. The exception is the FCoE application. If FCoE is the only application that you want the interface to advertise, then you do not need to define the FCoE application. You need to define the FCoE application only if you want interfaces to advertise other applications in addition to FCoE.

**NOTE:** If FCoE is the only application that you want DCBX to advertise on an interface, DCBX exchanges FCoE application protocol TLVs by default if the interface:

- Carries FCoE traffic (traffic mapped by CoS configuration to the FCoE forwarding class and applied to the interface)



- Has a congestion notification profile with PFC enabled on the FCoE priority (IEEE 802.1p code point)
- Does *not* have an application map

If you apply an application map to an interface, then all applications that you want DCBX to advertise must be defined and configured in the application map, including the FCoE application.

If no CoS configuration for FCoE is mapped to an interface, that interface does not exchange FCoE application protocol TLVs.

You can define:

- Layer 2 applications by EtherType
- Layer 4 applications by a combination of protocol (TCP or UDP) and destination port number

The EtherType is a two-octet field in the Ethernet frame that denotes the protocol encapsulated in the frame. For a list of common EtherTypes, see <http://standards.ieee.org/develop/regauth/ethertype/eth.txt> on the IEEE standards organization website. For a list of port numbers and protocols, see the *Service Name and Transport Protocol Port Number Registry* at [http://www.iana.org/assignments/service-names-port-numbers.xml](http://www.iana.org/assignments/service-names-port-numbers/service-names-port-numbers.xml) on the Internet Assigned Numbers Authority (IANA) website.

You must explicitly define each application that you want to advertise, except FCoE. The FCoE application is defined by default (EtherType 0x8906).

## Application Maps

An application map maps defined applications to one or more IEEE 802.1p code points. Each application map contains one or more applications. DCBX includes the configured application code points in the protocol TLVs exchanged with the connected peer.

To exchange protocol TLVs for an application, you must include the application in an application map. The FCoE application is a special case:

- If you want DCBX to exchange application protocol TLVs for more than one application on a particular interface, you must configure the applications, define an application map to map the applications to code points, and apply the application map to the interface. In this case, you must also define the FCoE application and add it to the application map.

This is the same process and treatment required for all other applications. In addition, for DCBX to exchange FCoE application TLVs, you must enable *priority-based flow control* (PFC) on the FCoE priority (the FCoE IEEE 802.1p code point) on the interface.

- If FCoE is the only application that you want DCBX to advertise on an interface, then you do not need to configure an application map and apply it to the interface. By default, when an interface has no application map, and the interface carries traffic mapped to the FCoE forwarding class, and PFC is enabled on the FCoE priority, the interface advertises FCoE TLVs (autonegotiation mode). DCBX exchanges FCoE application protocol TLVs by default until you apply an application map to the interface, remove the FCoE traffic from the interface (you can do this by removing the or editing the classifier for FCoE traffic), or disable PFC on the FCoE priority.

If you apply an application map to an interface that did not have an application map and was exchanging FCoE application TLVs, and you do not include the FCoE application in the application map, the interface stops exchanging FCoE TLVs. Every interface that has an application map must have FCoE included in the application map (and PFC enabled on the FCoE priority) in order for DCBX to exchange FCoE TLVs.

Mapping an application to code points does two things:

- Maps incoming traffic with the same code points to that application
- Allows you to configure classifiers that map incoming application traffic, by code point, to a forwarding class and a loss priority, in order to apply *class of service* (CoS) to application traffic and prioritize application traffic

You apply an application map to an interface to enable DCBX application protocol exchange on that interface for each application specified in the application map. All of the applications that you want an interface to advertise must be configured in the application map that you apply to the interface, with the previously noted exception for the FCoE application when FCoE is the only application for which you want DCBX to exchange protocol TLVs on an interface.

## Classifying and Prioritizing Application Traffic

When traffic arrives at an interface, the interface classifies the incoming traffic based on its code points. Classifiers map code points to loss priorities and forwarding classes. The loss priority prioritizes the traffic. The forwarding class determines the traffic output queue and CoS service level.

When you map an application to an IEEE 802.1p code point in an application map and apply the application map to an interface, incoming traffic on the interface that matches the application code points is mapped to the appropriate application. The application receives the loss priority and the CoS associated with the forwarding class for those code points, and is placed in the output queue associated with the forwarding class.

You can use the default classifier or you can configure a classifier to map the application code points defined in the application map to forwarding classes and loss priorities.

## Enabling Interfaces to Exchange Application Protocol Information

Each interface with the `fcoe` forwarding class and PFC enabled on the FCoE code point is enabled for FCoE application protocol exchange by default until you apply an application map to the interface. If you apply an application map to an interface and you want that interface to exchange FCoE application protocol TLVs, you must include the FCoE application in the application map. (In all cases, to achieve lossless transport, you must also enable PFC on the FCoE code point or code points.)

Except when FCoE is the only protocol you want DCBX to advertise on an interface, interfaces on which you want to exchange application protocol TLVs must include the following two items:

- The application map that contains the application(s)
- A classifier

**NOTE:** You must also enable PFC on the code point of any traffic for which you want to achieve lossless transport.

## Disabling DCBX Application Protocol Exchange

To disable DCBX application protocol exchange for all applications on an interface, issue the `set protocols dcbx interface interface-name applications no-auto-negotiation` command.

You can also disable DCBX application protocol exchange for applications on an interface by deleting the application map from the interface, or by deleting a particular application from the application map. However, when you delete an application from an application map, the application protocol is no longer exchanged on any interface which uses that application map.

On interfaces that use IEEE DCBX mode to exchange DCBX parameters, you can disable sending the enhanced transmission selection (ETS) Recommendation TLV to the peer if you want an asymmetric ETS configuration between the peers.

## RELATED DOCUMENTATION

*Understanding DCBX*

*Configuring DCBX Autonegotiation*

*Disabling the ETS Recommendation TLV*

*Defining an Application for DCBX Application Protocol TLV Exchange*

*Configuring an Application Map for DCBX Application Protocol TLV Exchange*

*Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange*

## Defining an Application for DCBX Application Protocol TLV Exchange

Define each application for which you want DCBX to exchange application protocol information. You can define Layer 2 and Layer 4 applications. After you define applications, you map them to IEEE 802.1p code points, and then apply the application map to the interfaces on which you want DCBX to exchange application protocol information with connected peers. (See *Related Documentation* for how to configure application maps and apply them to interfaces, and for an example of the entire procedure that also includes classifier configuration.)

**NOTE:** In Junos OS Release 12.1, the FCoE application was configured by default, so you did not need to configure it in an application map. In Junos OS Release 12.2, if you want DCBX to advertise the FCoE application on an interface and you apply an application map to that interface, you must explicitly configure FCoE in the application map. You also must enable priority-based flow control (PFC) on the FCoE code point on all interfaces that you want to advertise FCoE. If you apply an application map to an interface, the interface sends DCBX TLVs only for the applications configured in the application map.

Define Layer 2 applications by mapping an application name to an EtherType. Define Layer 4 applications by mapping an application name to a protocol (TCP or UDP) and a destination port.

- To define a Layer 2 application, specify the name of the application and its EtherType:

```
[edit applications]
user@switch# set application application-name ether-type ether-type
```

For example, to configure an application named PTP (for Precision Time Protocol) that uses the EtherType 0x88F7:

```
user@switch# set applications application ptp ether-type 0x88F7
```

- To define a Layer 4 application, specify the name of the application, its protocol (TCP or UDP), and its destination port:

```
[edit]
user@switch# set applications application application-name protocol (tcp | udp) destination-
port port-value
```

For example, to configure an application named `iscsi` (for Internet Small Computer System Interface) that uses the protocol TCP and the destination port 3260:

```
user@switch# set applications application iscsi protocol tcp destination-port 3260
```

## RELATED DOCUMENTATION

*Configuring an Application Map for DCBX Application Protocol TLV Exchange*

*Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange*

*Configuring DCBX Autonegotiation*

*Example: Configuring DCBX Application Protocol TLV Exchange*

[Example: Configuring DCBX to Support an iSCSI Application](#)

*Understanding DCBX Application Protocol TLV Exchange*

*show dcbx neighbors*

## Configuring an Application Map for DCBX Application Protocol TLV Exchange

After you define applications for which you want to exchange DCBX application protocol information, map the applications to IEEE 802.1p code points. The IEEE 802.1p code points identify incoming traffic and allow you to map that traffic to the desired application. You then apply the application map to the interfaces on which you want DCBX to exchange application protocol information with connected peers. (See *Related Documentation* for how to define applications and apply the application map to interfaces, and for an example of the entire procedure that also includes classifier configuration.)

**NOTE:** In Junos OS Release 12.1, the FCoE application was configured by default, so you did not need to configure it in an application map. In Junos OS Release 12.2, if you want DCBX to advertise the FCoE application on an interface and you apply an application map to that interface, you must explicitly configure FCoE in the application map. You also must enable priority-based flow control (PFC) on the FCoE code point on all interfaces that you want to advertise FCoE. If you apply an application map to an interface, the interface sends DCBX TLVs only for the applications configured in the application map.

Configure an application map by creating an application map name and mapping an application to one or more IEEE 802.1p code points.

- To define an application map, specify the name of the application map, the name of the application, and the IEEE 802.1p code points of the incoming traffic that you want to associate with the application in the application map:

```
[edit policy-options]
user@switch# set application-maps application-map-name application application-name code-
points [aliases] [bit-patterns]
```

For example, to configure an application map named ptp-app-map that includes an application named PTP (for Precision Time Protocol) and map the application to IEEE 802.1p code points 001 and 101:

```
user@switch# set policy-options application-maps ptp-app-map application ptp code points
[001 101]
```

## RELATED DOCUMENTATION

*Defining an Application for DCBX Application Protocol TLV Exchange*

*Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange*

*Configuring DCBX Autonegotiation*

*Example: Configuring DCBX Application Protocol TLV Exchange*

[Example: Configuring DCBX to Support an iSCSI Application](#)

*show dcbx neighbors*

## Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange

After you define applications and map them to IEEE 802.1p code points in an application map, apply the application map to the interfaces on which you want DCBX to exchange the application protocol information with connected peers. (See *Related Documentation* for how to define applications and configure application maps to interfaces, and for an example of the entire procedure that also includes classifier configuration.)

**NOTE:** In Junos OS Release 12.1, the FCoE application was configured by default, so you did not need to configure it in an application map. In Junos OS Release 12.2, if you want DCBX to advertise the FCoE application on an interface and you apply an application map to that interface, you must explicitly configure FCoE in the application map. You also must enable priority-based flow control (PFC) on the FCoE code point on all interfaces that you want to advertise FCoE. If you apply an application map to an interface, the interface sends DCBX TLVs only for the applications configured in the application map.

- To apply an application map to a DCBX interface, specify the DCBX interface and the application map name:

```
[edit protocols]
user@switch# set dcbx interface interface-name application-map application-map-name
```

For example, to apply an application map named ptp-app-map on interface xe-0/0/11:

```
user@switch# set protocols dcbx interface xe-0/0/11 application-map ptp-app-map
```

### RELATED DOCUMENTATION

*Defining an Application for DCBX Application Protocol TLV Exchange*

*Configuring an Application Map for DCBX Application Protocol TLV Exchange*

*Configuring DCBX Autonegotiation*

*Example: Configuring DCBX Application Protocol TLV Exchange*

[Example: Configuring DCBX to Support an iSCSI Application](#)

*show dcbx neighbors*

## Example: Configuring DCBX Application Protocol TLV Exchange

### IN THIS SECTION

- [Requirements | 662](#)
- [Overview | 662](#)
- [Configuration | 667](#)
- [Verification | 670](#)

Data Center Bridging Capability Exchange protocol (DCBX) discovers the data center bridging (DCB) capabilities of connected peers by exchanging application configuration information. DCBX detects feature misconfiguration and mismatches and can configure DCB on peers. DCBX is an extension of the Link Layer Discovery Protocol (LLDP). LLDP must remain enabled on every interface on which you want to use DCBX.

**NOTE:** LLDP and DCBX are enabled by default on all interfaces.

The switch supports DCBX application protocol exchange for Layer 2 and Layer 4 applications such as the Internet Small Computer System Interface (iSCSI). You specify applications by EtherType (for Layer 2 applications) or by the destination port and protocol (for Layer 4 applications; the protocol can be either TCP or UDP).

The switch handles Fibre Channel over Ethernet (FCoE) application protocol exchange differently than other protocols in some cases:

- If FCoE is the only application for which you want to enable DCBX application protocol TLV exchange on an interface, you do not have to explicitly configure the FCoE application or an application map. By default, the switch exchanges FCoE application protocol TLVs on all interfaces that carry FCoE traffic (traffic mapped to the `fcoe` forwarding class) and have priority-based flow control (PFC) enabled on the FCoE priority (the FCoE IEEE 802.1p code point). The default priority mapping for the FCoE application is IEEE 802.1p code point 011 (the default `fcoe` forwarding class code point).
- If you want an interface to use DCBX to exchange application protocol TLVs for any other applications in addition to FCoE, you must configure the applications (including FCoE), define an application map (including FCoE), and apply the application map to the interface. If you apply an application map to an interface, you must explicitly configure the FCoE application, or the interface does not exchange FCoE application protocol TLVs.



This example shows how to configure interfaces to exchange both Layer 2 and Layer 4 applications by configuring one interface to exchange iSCSI and FCoE application protocol information and configuring another interface to exchange iSCSI and Precision Time Protocol (PTP) application protocol information.

## Requirements

This example uses the following hardware and software components:

- Juniper Networks QFX Series device
- Junos OS Release 12.1 or later for the QFX Series

## Overview

### IN THIS SECTION

- [Topology](#) | 664

The switch supports DCBX application protocol exchange for:

- Layer 2 applications, defined by EtherType
- Layer 4 applications, defined by destination port and protocol

**NOTE:** DCBX also advertises PFC and enhanced transmission selection (ETS) information. See [Configuring DCBX Autonegotiation](#) for how DCBX negotiates and advertises configuration information for these features and for the applications.

DCBX is configured on a per-interface basis for each supported feature or application. For applications that you want to enable for DCBX application protocol exchange, you must:

- Define the application name and configure the EtherType or the destination port and protocol (TCP or UDP) of the application. Use the EtherType for Layer 2 applications, and use the destination port and protocol for Layer 4 protocols.
- Map the application to an IEEE 802.1p code point in an application map.
- Add the application map to DCBX interface.

In addition, for all applications (including FCoE, even when you do not use an application map), you either must create an IEEE 802.1p classifier and apply it to the appropriate ingress interfaces or use the default classifier. A classifier maps the code points of incoming traffic to a forwarding class and a loss

priority so that ingress traffic is assigned to the correct class of service (CoS). The forwarding class determines the output queue on the egress interface.

If you do not create classifiers, trunk and tagged-access ports use the unicast IEEE 802.1 default trusted classifier. [Table 102 on page 663](#) shows the default mapping of IEEE 802.1 code-point values to unicast forwarding classes and loss priorities for ports in trunk mode or tagged-access mode. [Table 103 on page 663](#) shows the default untrusted classifier IEEE 802.1 code-point values to unicast forwarding class mapping for ports in access mode.

**Table 102: Default IEEE 802.1 Classifiers for Trunk Ports and Tagged-Access Ports (Default Trusted Classifier)**

| Code Point | Forwarding Class | Loss Priority |
|------------|------------------|---------------|
| be (000)   | best-effort      | low           |
| be1 (001)  | best-effort      | low           |
| ef (010)   | best-effort      | low           |
| ef1 (011)  | fcoe             | low           |
| af11 (100) | no-loss          | low           |
| af12 (101) | best-effort      | low           |
| nc1 (110)  | network-control  | low           |
| nc2 (111)  | network-control  | low           |

**Table 103: Default IEEE 802.1 Unicast Classifiers for Access Ports (Default Untrusted Classifier)**

| Code Point | Forwarding Class | Loss Priority |
|------------|------------------|---------------|
| 000        | best-effort      | low           |
| 001        | best-effort      | low           |

**Table 103: Default IEEE 802.1 Unicast Classifiers for Access Ports (Default Untrusted Classifier)**  
(Continued)

| Code Point | Forwarding Class | Loss Priority |
|------------|------------------|---------------|
| 010        | best-effort      | low           |
| 011        | best-effort      | low           |
| 100        | best-effort      | low           |
| 101        | best-effort      | low           |
| 110        | best-effort      | low           |
| 111        | best-effort      | low           |

### Topology

This example shows how to configure DCBX application protocol exchange for three protocols (iSCSI, PTP, and FCoE) on two interfaces. One interface exchanges iSCSI and FCoE application protocol information, and the other interface exchanges iSCSI and PTP application protocol information.

**NOTE:** You must map FCoE traffic to the interfaces on which you want to forward FCoE traffic. You must also enable PFC on the FCoE interfaces and create an ingress classifier for FCoE traffic, or else use the default classifier.

[Table 104 on page 664](#) shows the configuration components for this example.

**Table 104: Components of DCBX Application Protocol Exchange Configuration Topology**

| Component | Settings          |
|-----------|-------------------|
| Hardware  | QFX Series device |

**Table 104: Components of DCBX Application Protocol Exchange Configuration Topology** *(Continued)*

| Component                   | Settings                                                                                                                                                                                                                                                                                                                             |
|-----------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| LLDP                        | Enabled by default on Ethernet interfaces                                                                                                                                                                                                                                                                                            |
| DCBX                        | Enabled by default on Ethernet interfaces                                                                                                                                                                                                                                                                                            |
| iSCSI application (Layer 4) | Application name—iscsi<br>protocol—TCP<br>destination-port—3260<br>code-points—111                                                                                                                                                                                                                                                   |
| PTP application (Layer 2)   | Application name—ptp<br>ether-type—0x88F7<br>code-points—001, 101                                                                                                                                                                                                                                                                    |
| FCoE application (Layer 2)  | Application name—fcoe<br>ether-type—0x8906<br>code-points—011<br><br><b>NOTE:</b> You explicitly configure the FCoE application because you are applying an application map to the interface. When you apply an application map to an interface, all applications must be explicitly configured and included in the application map. |
| Application maps            | dcbx-iscsi-fcoe-app-map—Maps the iSCSI and FCoE applications to IEEE 802.1p code points<br><br>dcbx-iscsi-ptp-app-map—Maps iSCSI and PTP applications to IEEE 802.1p code points                                                                                                                                                     |

**Table 104: Components of DCBX Application Protocol Exchange Configuration Topology** *(Continued)*

| Component                                                                                                         | Settings                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
|-------------------------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Interfaces                                                                                                        | <p>xe-0/0/10—Configured to exchange FCoE and iSCSI application TLVs (uses application map dcbx-iscsi-fcoe-app-map, carries FCoE traffic, and has PFC enabled on the FCoE priority)</p> <p>xe-0/0/11—Configured to exchange iSCSI and PTP application TLVs (uses application map dcbx-iscsi-ptp-app-map)</p>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                             |
| PFC congestion notification profile for FCoE application exchange                                                 | <p>fcoe-cnp:</p> <ul style="list-style-type: none"> <li>• Code point—011</li> <li>• Interface—xe-0/0/10</li> </ul>                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                      |
| Behavior aggregate classifiers (map forwarding classes to incoming packets by the packet's IEEE 802.1 code point) | <p>fcoe-iscsi-cl1:</p> <ul style="list-style-type: none"> <li>• Maps the fcoe forwarding class to the IEEE 802.1p code point used for the FCoE application (011) and a loss priority of high</li> <li>• Maps the network-control forwarding class to the IEEE 802.1p code point used for the iSCSI application (111) and a loss priority of high</li> <li>• Applied to interface xe-0/0/10</li> </ul> <p>iscsi-ptp-cl2:</p> <ul style="list-style-type: none"> <li>• Maps the network-control forwarding class to the IEEE 802.1p code point used for the iSCSI application (111) and a loss priority of low</li> <li>• Maps the best-effort forwarding class to the IEEE 802.1p code points used for the PTP application (001 and 101) and a loss priority of low</li> <li>• Applied to interface xe-0/0/11</li> </ul> |

**NOTE:** This example does not include scheduling (bandwidth allocation) configuration or lossless configuration for the iSCSI forwarding class.

## Configuration

### IN THIS SECTION

- [CLI Quick Configuration | 667](#)
- [Configuring DCBX Application Protocol TLV Exchange | 668](#)

### CLI Quick Configuration

To quickly configure DCBX application protocol exchange, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level.

```
set applications application iSCSI protocol tcp destination-port 3260
set applications application FCoE ether-type 0x8906
set applications application PTP ether-type 0x88F7
set policy-options application-maps dcbx-iscsi-fcoe-app-map application iSCSI code-points 111
set policy-options application-maps dcbx-iscsi-fcoe-app-map application FCoE code-points 011
set policy-options application-maps dcbx-iscsi-ptp-app-map application iSCSI code-points 111
set policy-options application-maps dcbx-iscsi-ptp-app-map application PTP code-points [001 101]
set protocols dcbx interface xe-0/0/10 application-map dcbx-iscsi-fcoe-app-map
set protocols dcbx interface xe-0/0/11 application-map dcbx-iscsi-ptp-app-map
set class-of-service congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
set class-of-service interfaces xe-0/0/10 congestion-notification-profile fcoe-cnp
set class-of-service classifiers ieee-802.1 fcoe-iscsi-cl1 import default forwarding-class fcoe
loss-priority high code-points 011
set class-of-service classifiers ieee-802.1 fcoe-iscsi-cl1 import default forwarding-class
network-control loss-priority high code-points 111
set class-of-service classifiers ieee-802.1 iscsi-ptp-cl2 import default forwarding-class
network-control loss-priority low code-points 111
set class-of-service classifiers ieee-802.1 iscsi-ptp-cl2 import default forwarding-class best-
effort loss-priority low code-points [001 101]
set class-of-service interfaces xe-0/0/10 unit 0 classifiers ieee-802.1 fcoe-iscsi-cl1
set class-of-service interfaces xe-0/0/11 unit 0 classifiers ieee-802.1 iscsi-ptp-cl2
```

## Configuring DCBX Application Protocol TLV Exchange

### Step-by-Step Procedure

To define the applications, map the applications to IEEE 802.1p code points, apply the applications to interfaces, and create classifiers for DCBX application protocol exchange:

1. Define the iSCSI application by specifying its protocol and destination port, and define the FCoE and PTP applications by specifying their EtherTypes.

```
[edit applications]
user@switch# set application iSCSI protocol tcp destination-port 3260
user@switch# set application FCoE ether-type 0x8906
user@switch# set application PTP ether-type 0x88F7
```

2. Define an application map that maps the iSCSI and FCoE applications to IEEE 802.1p code points.

```
[edit policy-options]
user@switch# set application-maps dcbx-iscsi-fcoe-app-map application iSCSI code-points 111
user@switch# set application-maps dcbx-iscsi-fcoe-app-map application FCoE code-points 011
```

3. Define the application map that maps the iSCSI and PTP applications to IEEE 802.1p code points.

```
[edit policy-options]
user@switch# set application-maps dcbx-iscsi-ntp-app-map application iSCSI code-points 111
user@switch# set application-maps dcbx-iscsi-ntp-app-map application PTP code-points [001 101]
```

4. Apply the iSCSI and FCoE application map to interface xe-0/0/10, and apply the iSCSI and PTP application map to interface xe-0/0/11.

```
[edit protocols dcbx]
user@switch# set interface xe-0/0/10 application-map dcbx-iscsi-fcoe-app-map
user@switch# set interface xe-0/0/11 application-map dcbx-iscsi-ntp-app-map
```

5. Create the congestion notification profile to enable PFC on the FCoE code point (011), and apply the congestion notification profile to interface xe-0/0/10.

```
[edit class-of-service]
user@switch# set congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
user@switch# set interfaces xe-0/0/10 congestion-notification-profile fcoe-cnp
```

6. Configure the classifier to apply to the interface that exchanges iSCSI and FCoE application information.

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 fcoe-iscsi-cl1 import default forwarding-class fcoe loss-priority
high code-points 011
user@switch# set ieee-802.1 fcoe-iscsi-cl1 import default forwarding-class network-control
loss-priority high code-points 111
```

7. Configure the classifier to apply to the interface that exchanges iSCSI and PTP application information.

```
[edit class-of-service classifiers]
user@switch# set ieee-802.1 iscsi-ntp-cl2 import default forwarding-class network-control
loss-priority low code-points 111
user@switch# set ieee-802.1 iscsi-ntp-cl2 import default forwarding-class best-effort loss-
priority low code-points [001 101]
```

8. Apply the classifiers to the appropriate interfaces.

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/10 unit 0 classifiers ieee-802.1 fcoe-iscsi-cl1
user@switch# set interfaces xe-0/0/11 unit 0 classifiers ieee-802.1 iscsi-ntp-cl2
```



## Verification

### IN THIS SECTION

- [Verifying the Application Configuration | 670](#)
- [Verifying the Application Map Configuration | 671](#)
- [Verifying DCBX Application Protocol Exchange Interface Configuration | 672](#)
- [Verifying the PFC Configuration | 672](#)
- [Verifying the Classifier Configuration | 673](#)

To verify that DCBX application protocol exchange configuration has been created and is operating properly, perform these tasks:

### Verifying the Application Configuration

#### Purpose

Verify that DCBX applications have been configured.

#### Action

List the applications by using the configuration mode command `show applications`:

```
user@switch# show applications
application iSCSI {
 protocol tcp;
 destination-port 3260;
}

application fcoe {
 ether-type 0x8906;
}

application ptp {
 ether-type 0x88F7;
}
```

## Meaning

The `show applications configuration mode` command lists all of the configured applications and either their protocol and destination port (Layer 4 applications) or their EtherType (Layer 2 applications). The command output shows that the iSCSI application is configured with the `tcp` protocol and destination port 3260, the FCoE application is configured with the EtherType 0x8906, and that the PTP application is configured with the EtherType 0x88F7.

## Verifying the Application Map Configuration

### Purpose

Verify that the application maps have been configured.

### Action

List the application maps by using the configuration mode command `show policy-options application-maps`:

```
user@switch# show policy-options application-maps
dcbx-iscsi-fcoe-app-map {
 application iSCSI code-points 111;
 application FCoE code-points 011;
}

dcbx-iscsi-ptp-app-map {
 application iSCSI code-points 111;
 application PTP code-points [001 101];
}
```

## Meaning

The `show policy-options application-maps` configuration mode command lists all of the configured application maps and the applications that belong to each application map. The command output shows that there are two application maps, `dcbx-iscsi-fcoe-app-map` and `dcbx-iscsi-ptp-app-map`.

The application map `dcbx-iscsi-fcoe-app-map` consists of the iSCSI application, which is mapped to IEEE 802.1p code point 111, and the FCoE application, which is mapped to IEEE 802.1p code point 011.

The application map `dcbx-iscsi-ptp-app-map` consists of the iSCSI application, which is mapped to IEEE 802.1p code point 111, and the PTP application, which is mapped to IEEE 802.1p code points 001 and 101.

## Verifying DCBX Application Protocol Exchange Interface Configuration

### Purpose

Verify that the application maps have been applied to the correct interfaces.

### Action

List the application maps by using the configuration mode command `show protocols dcbx`:

```
user@switch# show protocols dcbx
interface xe-0/0/10.0 {
 application-map dcbx-iscsi-fcoe-app-map;
}

interface xe-0/0/11.0 {
 application-map dcbx-iscsi-ptp-app-map;
}
```

### Meaning

The `show protocols dcbx` configuration mode command lists whether the interfaces are enabled for DCBX and lists the application map applied to each interface. The command output shows that interfaces `xe-0/0/10.0` and `xe-0/0/11.0` are enabled for DCBX, and that interface `xe-0/0/10.0` uses application map `dcbx-iscsi-fcoe-app-map`, and interface `xe-0/0/11.0` uses application map `dcbx-iscsi-ptp-app-map`.

## Verifying the PFC Configuration

### Purpose

Verify that PFC has been enabled on the FCoE code point and applied to the correct interface.

### Action

Display the PFC configuration to verify that PFC is enabled on the FCoE code point (011) in the congestion notification profile `fcoe-cnp` by using the configuration mode command `show class-of-service congestion-notification-profile`:

```
user@switch# show class-of-service congestion-notification-profile
fcoe-cnp {
```

```

input {
 ieee-802.1 {
 code-point 011 {
 pfc;
 }
 }
}

```

Display the class-of-service (CoS) interface information to verify that the correct interface has PFC enabled for the FCoE application by using the configuration mode command `show class-of-service interfaces`:

```

user@switch# show class-of-service interfaces
xe-0/0/10 {
 congestion-notification-profile fcoe-cnp;
}

```

**NOTE:** The sample output does not include all of the information this command can show. The output is abbreviated to focus on verifying the PFC configuration.

## Meaning

The `show class-of-service congestion-notification-profile` configuration mode command lists the configured congestion notification profiles. The command output shows that the congestion notification profile `fcoe-cnp` has been configured and has enabled PFC on the IEEE 802.1p code point 011 (the default FCoE code point).

The `show class-of-service interfaces` configuration mode command shows the interface CoS configuration. The command output shows that the congestion notification profile `fcoe-cnp`, which enables PFC on the FCoE code point, is applied to interface `xe-0/0/10`.

## Verifying the Classifier Configuration

### Purpose

Verify that the classifiers have been configured and applied to the correct interfaces.

## Action

Display the classifier configuration by using the configuration mode command `show class-of-service`:

```
user@switch# show class-of-service
classifiers {
 ieee-802.1 fcoe-iscsi-cl1 {
 import default;
 forwarding-class network-control {
 loss-priority high code-points 111;
 }
 forwarding-class fcoe {
 loss-priority high code-points 011;
 }
 }
 ieee-802.1 iscsi-ntp-cl2 {
 import default;
 forwarding-class network-control {
 loss-priority low code-points 111;
 }
 forwarding-class best-effort {
 loss-priority low code-points [001 101];
 }
 }
}
interfaces {
 xe-0/0/10 {
 congestion-notification-profile fcoe-cnp;
 unit 0 {
 classifiers {
 ieee-802.1 fcoe-iscsi-cl1;
 }
 }
 }
 xe-0/0/11 {
 unit 0 {
 classifiers {
 ieee-802.1 iscsi-ntp-cl2;
 }
 }
 }
}
```

```
}
}
```

**NOTE:** The sample output does not include all of the information this command can show. The output is abbreviated to focus on verifying the classifier configuration.

## Meaning

The `show class-of-service configuration mode` command lists the classifier and CoS interface configuration, as well as other information not shown in this example. The command output shows that there are two classifiers configured, `fcoe-iscsi-cl1` and `iscsi-ntp-cl2`.

Classifier `fcoe-iscsi-cl1` uses the default classifier as a template and edits the template as follows:

- The forwarding class `network-control` is set to a loss priority of `high` and is mapped to code point 111 (the code point mapped to the iSCSI application).
- The forwarding class `fcoe` is set to a loss priority of `high` and is mapped to code point 011 (the code point mapped by default to the FCoE application).

Classifier `iscsi-ntp-cl2` uses the default classifier as a template and edits the template as follows:

- The forwarding class `network-control` is set to a loss priority of `low` and is mapped to IEEE 802.1p code point 111 (the code point mapped to the iSCSI application).
- The forwarding class `best-effort` is set to a loss priority of `low` and is mapped to IEEE 802.1p code points 001 and 101 (the code points mapped by default to the PTP application).

The command output also shows that classifier `fcoe-iscsi-cl1` is mapped to interface `xe-0/0/10.0` and that classifier `iscsi-ntp-cl2` is mapped to interface `xe-0/0/11.0`.

## RELATED DOCUMENTATION

*Defining an Application for DCBX Application Protocol TLV Exchange*

*Configuring an Application Map for DCBX Application Protocol TLV Exchange*

*Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange*

*Configuring DCBX Autonegotiation*

*show dcbx*

*show dcbx neighbors*

*Understanding DCBX Application Protocol TLV Exchange*

# 5

PART

## Configuring Buffers

---

[Using Buffers](#) | 677

---

## CHAPTER 8

# Using Buffers

**IN THIS CHAPTER**

- [Understanding CoS Buffer Configuration | 677](#)
- [Configuring Global Ingress and Egress Shared Buffers | 701](#)
- [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic | 703](#)
- [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled | 712](#)
- [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 720](#)
- [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 729](#)

## Understanding CoS Buffer Configuration

**IN THIS SECTION**

- [Buffer Pools | 679](#)
- [Default Buffer Pool Values | 689](#)
- [Shared Buffer Configuration Recommendations for Different Network Traffic Scenarios | 693](#)
- [Optimizing Buffer Configuration | 698](#)
- [General Buffer Configuration Rules and Considerations | 699](#)

Packet Forwarding Engine (PFE) wide common packet buffer memory is used to store packets on interface queues. The buffer memory has separate ingress and egress accounting to make accept, drop, or pause decisions. Because the switch has a single pool of memory with separate ingress and egress accounting, the full amount of buffer memory is available from both the ingress and the egress



perspective. Packets are accounted for as they enter and leave the switch, but there is no concept of a packet arriving at an ingress buffer and then being moved to an egress buffer. Specific common buffer memory amounts for individual switches is listed in [Table 105 on page 678](#).

**Table 105: Common Packet Buffer Memory on Switches**

| Switch                          | Common Packet Buffer Memory |
|---------------------------------|-----------------------------|
| QFX3500, QFX3600                | 9MB                         |
| QFX5100, EX4600, and OCX Series | 12MB                        |
| QFX5110, QFX5200-32C            | 16MB                        |
| QFX5200-48Y                     | 22MB                        |
| QFX5120                         | 32MB                        |
| QFX5130, QFX5700                | 132MB                       |
| QFX5210                         | 42MB                        |
| QFX5220                         | 64MB                        |

**NOTE:** QFX10000 does not have a shared buffer.

The buffers are divided into two pools from both an ingress and an egress perspective:

1. *Shared buffers* are a global memory pool that the switch allocates dynamically to ports as needed, so the buffers are shared among the switch ports.
2. *Dedicated buffers* are a memory pool divided equally among the switch ports. Each port receives a minimum guaranteed amount of buffer space, dedicated to each port, not shared among ports.

**NOTE:** Lossless traffic is traffic on which you enable priority-based flow control (PFC) to ensure lossless transport. Lossless traffic does not refer to best-effort traffic on a link enabled for Ethernet PAUSE (IEEE 802.3x).

The switch reserves nonconfigurable buffer space to ensure that ports and queues receive a minimum memory allocation. You can configure how the system uses the rest of the buffer space to optimize the allocation for your mix of network traffic. You can configure the percentage of available buffer space used as shared buffer space versus dedicated buffer space. You can also configure how shared buffer space is allocated to different types of traffic. You can optimize the buffer settings for the traffic on your network.

The default class-of-service configuration provides two lossless forwarding classes (`fcoe` and `no-loss`), a best-effort unicast forwarding class, a network control traffic forwarding class, and one multidestination (multicast, broadcast, and destination lookup fail) forwarding class.

Each default forwarding class maps to a different default output queue. The default configuration allocates the buffers in a manner that supports a moderate amount of lossless traffic while still providing the ability to absorb bursts in best-effort traffic transmission.

Changing the buffer settings changes the abilities of the buffers to absorb traffic bursts and handle lossless traffic. For example, networks with mostly best-effort traffic require allocating most of the shared buffer space to best-effort buffers. This provides deep, flexible buffers that can absorb traffic bursts with minimal packet loss, at the expense of buffer availability for lossless traffic.

Conversely, networks with mostly lossless traffic require allocating most of the shared buffer space to lossless headroom buffers. This prevents packet loss on lossless flows at the expense of absorbing bursty best-effort traffic efficiently.



**CAUTION:** Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete.

This topic describes the buffer architecture and settings:

## Buffer Pools

From both an ingress and an egress perspective, the PFE buffer is split into two main pools, a shared buffer pool and a dedicated buffer pool that ensures a minimum allocation to each port. You can configure the amount of buffer space allocated to each of the two pools. A portion of the buffer space is reserved so that there is always a minimum amount of shared and dedicated buffer space available to each port.

- **Shared buffer pool**—A global memory space that all of the ports on the switch share dynamically as they need buffers. The shared buffer pool is further partitioned into buffers for best-effort unicast, best-effort multdestination (broadcast, multicast, and destination lookup fail), and PFC (lossless) traffic types. You can allocate global shared memory space to buffer partitions to better support different mixes of network traffic. The larger the shared buffer pool, the better the switch can absorb traffic bursts because more shared memory is available for the traffic.
- **Dedicated buffer pool**—A reserved global memory space allocated equally to each port. The switch reserves a minimum dedicated buffer pool that is not user-configurable. You can divide the dedicated buffer allocation for a port among the port queues on a per-port, per-queue basis. (For example, this enables you to dedicate more buffer space to queues that transport lossless traffic.)

A larger dedicated buffer pool means a larger amount of dedicated buffer space for each port, so congestion on one port is less likely to affect traffic on another port because the traffic does not need to use as much shared buffer space. However, the larger the dedicated buffer pool, the less bursty traffic the switch can handle because there is less dynamic shared buffer memory.

You can configure the way the available unreserved portion of the buffer space is allocated to the global shared buffer pool and to the dedicated shared buffer pool by configuring the ingress and egress shared buffer percentages.

By default, 100 percent of the available unreserved buffer space is allocated to the shared buffer pool. If you change the percentage of space allocated to the shared buffer, the available buffer space that is not allocated to the shared buffer is allocated to the dedicated buffer. For example, if you configure the ingress shared buffer pool as 80 percent, the remaining 20 percent of the available buffer space is allocated to the dedicated buffer pool and divided equally across the ports.

**NOTE:** When 100 percent of the available (user-configurable) buffers are allocated to the shared buffer pool, the switch still reserves a minimum dedicated buffer pool.

You can separately configure ingress and egress shared buffer pool allocations. You can also partition the ingress and egress shared buffer pool to allocate percentages of the shared buffer pool to specific types of traffic. If you do not use the default configuration or one of the recommended configurations, pay particular attention to the ingress configuration of the lossless headroom buffers (these buffers handle PFC pause during periods of congestion) and to the egress configuration of the best-effort buffers to handle incast congestion (multiple synchronized sources sending data to the same receiver in parallel).

In addition to the shared buffer pool and the dedicated buffer pool, there is also a small ingress global headroom buffer pool that is reserved and is not configurable.

When contention for buffer space occurs, the switch uses an internal algorithm to ensure that the buffer pools are distributed fairly among competing flows. When traffic for a given flow exceeds the amount of dedicated port buffer reserved for that flow, the flow begins to consume memory from the dynamic

shared buffer pool. Competing flows compete for shared buffer memory with other flows that also have exhausted their dedicated buffers. When there is no congestion, there are no competing flows.

## Buffer Handling of Lossless Flows (PFC) Versus Ethernet PAUSE

When we discuss lossless buffers in the following sections, we mean buffers that handle traffic on which you enable PFC to ensure lossless transport. The lossless buffers are not used for best-effort traffic on a link on which you enable Ethernet PAUSE (IEEE 802.3x). The lossless ingress and egress shared buffers, and the ingress lossless headroom shared buffer, are used only for traffic on which you enable PFC.

**NOTE:** To support lossless flows, you must configure the appropriate data center bridging capabilities (PFC, DCBX, and ETS) and scheduling properties.

## Shared Buffer Pool and Partitions

The shared buffer pool is a global memory space that all of the ports on the switch share dynamically as they need buffers. The switch uses the shared buffer pool to absorb traffic bursts after the dedicated buffer pool for a port is exhausted.

You can divide both the ingress shared buffer pool and the egress shared buffer pool into three partitions to allocate percentages of each buffer pool to different types of traffic. When you partition the ingress or egress shared buffer pool:

- If you explicitly configure one ingress shared buffer partition, you must explicitly configure all three ingress shared buffer partitions. (You either explicitly configure all three ingress partitions or you use the default setting for all three ingress partitions.)

If you explicitly configure one egress shared buffer partition, you must explicitly configure all three egress shared buffer partitions. (You either explicitly configure all three egress partitions or you use the default setting for all three egress partitions.)

The switch returns a commit error if you do not explicitly configure all three partitions when configuring the ingress or egress shared buffer partitions.

- The combined percentages of the three ingress shared buffer partitions must total exactly 100 percent.

The combined percentages of the three egress shared buffer partitions must total exactly 100 percent.

When you explicitly configure ingress or egress shared buffer partitions, the switch returns a commit error if the total percentage of the three partitions does not equal 100 percent.

- If you explicitly partition one set of shared buffers, you do not have to explicitly partition the other set of shared buffers. For example, you can explicitly configure the ingress shared buffer partitions and use the default egress shared buffer partitions. However, if you change the buffer partitions for the ingress buffer pool to match the expected types of traffic flows, you would probably also want to change the buffer partitions for the egress buffer pool to match those traffic flows.

You can configure the percentage of available unreserved buffer space allocated to the shared buffer pool. Space that you do not allocate to the shared buffer pool is added to the dedicated buffer pool and divided equally among the ports. The default configuration allocates 100 percent of the unreserved ingress and egress buffer space to the shared buffers.

Configuring the ingress and egress shared buffer pool partitions enables you to allocate more buffers to the types of traffic your network predominantly carries, and fewer buffers to other traffic.

### Ingress Shared Buffer Pool Partitions

You can configure three ingress buffer pool partitions:

- Lossless buffers—Shared buffer pool for all lossless ingress traffic. We recommend 5 percent as the minimum value for lossless buffers.
- Lossless headroom buffers—Shared buffer pool for packets received while a pause is asserted. If PFC is enabled on priorities on a port, when the port sends a pause message to the connected peer, the port uses the headroom buffers to store the packets that arrive between the time the port sends the pause message and the time the last packet arrives after the peer pauses traffic. The minimum value for lossless headroom buffers is 0 (zero) percent. (Lossless headroom buffers are the only buffers for which the recommended value can be less than 5 percent.)

**NOTE:** On a QFX Virtual Chassis and an EX4600/EX4650 Virtual Chassis, the minimum value for the lossless headroom buffer is 3 percent.

- Lossy buffers—Shared buffer pool for all best-effort ingress traffic (best-effort unicast, multidestination, and strict-high priority traffic). We recommend 5 percent as the minimum value for best-effort buffers.

The combined percentage values of the ingress lossless, lossless headroom, and best-effort buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. If you explicitly configure an ingress shared buffer partition, you must explicitly configure all three ingress buffer partitions, even if the lossless headroom buffer partition has a value of 0 (zero) percent.

### Egress Shared Buffer Pool Partitions

You can configure three egress buffer pool partitions:

- Lossless buffers—Shared buffer pool for all lossless egress queues. We recommend 5 percent as the minimum value for lossless buffers.
- Lossy buffers—Shared buffer pool for all best-effort egress queues (best-effort unicast, and strict-high priority queues). We recommend 5 percent as the minimum value for best-effort buffers.
- Multicast buffers—Shared buffer pool for all multidestination (multicast, broadcast, and destination lookup fail) egress queues. We recommend 5 percent as the minimum value for multicast buffers.

The combined percentage values of the egress lossless, lossy, and multicast buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All egress buffer partitions must be explicitly configured and should have a value of at least 5 percent. If you explicitly configure an egress shared buffer partition, you must explicitly configure all three egress buffer partitions, and each partition should have a value of at least 5 percent.

**NOTE:** QFX5200-32C does not replicate all multicast streams when two or more downstream interface packet sizes are higher than ~6k and have an 1000pps packet ingress rate. This is because the number of working flows on QFX5200-32C is indirectly proportional to the packet size and directly proportional to available multicast shared buffers.

## Dedicated Port Buffer Pool and Buffer Allocation to Queues

The global dedicated buffer pool is memory that is allocated equally to each port, so each port receives a guaranteed minimum amount of buffer space. Dedicated buffers are not shared among ports. Each port receives an equal proportion of the dedicated buffer pool.

When traffic enters and exits the switch, the switch ports use their dedicated buffers to store packets. If the dedicated buffers are not sufficient to handle the traffic, the switch uses shared buffers. The only way to increase the dedicated buffer pool is to decrease the shared buffer pool from its default value of 100 percent of available unreserved buffers.

The amount of dedicated buffer space is not user-configurable and depends on the percentage of available nonreserved buffers allocated to the shared buffers. (The dedicated buffer space is equal to the minimum reserved port buffers plus the remainder of the available nonreserved buffers that are not allocated to the shared buffer pool.)

**NOTE:** If 100 percent of the available unreserved buffers are allocated to the shared buffer pool, the switch still reserves a minimum dedicated buffer pool.

The larger the shared buffer pool, the better the burst absorption across the ports. The larger the dedicated buffer pool, the larger the amount of dedicated buffer space for each port. The greater the dedicated buffer space, the less likely that congestion on one port can affect traffic on another port, because the traffic does not need to use as much shared buffer space.

### Allocating Dedicated Port Buffers to Queues

You can divide the dedicated buffer allocation for an egress port among the port queues by including the `buffer-size` statement in the scheduler configuration. This enables you to control the egress port dedicated buffer allocation on a per-port, per-queue basis. (For example, this enables you to dedicate more buffer space to queues that transport lossless traffic, or to stop the port from reserving buffers for queues that do not carry traffic.) Egress dedicated port buffer allocation is a hierarchical structure that allocates a global dedicated buffer pool evenly among ports, and then divides the allocation for each port among the port queues.

By default, ports divide their allocation of dedicated buffers among their egress queues in the same proportion as the default scheduler sets the minimum guaranteed transmission rates (the `transmit-rate` option) for traffic. Only the queues included in the default scheduler receive bandwidth and dedicated buffers, in the proportions shown in [Table 106 on page 684](#):

**Table 106: Default Dedicated Buffer Allocation to Egress Queues (Based on Default Scheduler)**

| Forwarding Class | Queue | Minimum Guaranteed Bandwidth<br>( <code>transmit-rate</code> ) | Proportion of Reserved Dedicated Port Buffers |
|------------------|-------|----------------------------------------------------------------|-----------------------------------------------|
| best-effort      | 0     | 5%                                                             | 5%                                            |
| fcoe             | 3     | 35%                                                            | 35%                                           |
| no-loss          | 4     | 35%                                                            | 35%                                           |
| network-control  | 7     | 5%                                                             | 5%                                            |
| mcast            | 8     | 20%                                                            | 20%                                           |

In the default configuration, no egress queues other than the ones shown in [Table 106 on page 684](#) receive an allocation of dedicated port buffers.

**NOTE:** The switch uses hierarchical scheduling to control port and queue bandwidth allocation, as described in [Understanding CoS Hierarchical Port Scheduling \(ETS\)](#) and shown in [Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#). For egress queue buffer size configuration, when you attach a traffic control profile (includes the queue scheduler information) to a port, the dedicated egress buffers on the port are divided among the queues as configured in the scheduler.

If you do not want to use the default allocation of dedicated port buffers to queues, use the `buffer-size` option in the scheduler that is attached to the port to configure the queue allocation. You can configure the dedicated buffer allocation to queues in two ways:

- As a percentage—The queue receives the specified percentage of dedicated port buffers when the queue is mapped to the scheduler and the scheduler is attached to a port.
- As a remainder—After the port services the queues that have an explicit percentage buffer size configuration, the remaining dedicated port buffer space is divided equally among the other queues to which a scheduler is attached. (No default or explicit scheduler for a queue means no dedicated buffer allocation for that queue.) If you configure a scheduler and you do not specify a buffer size as a percentage, *remainder* is the default setting.

**NOTE:** The total of all of the explicitly configured buffer size percentages for all of the queues on a port cannot exceed 100 percent.

On all QFX5000 platforms, when calculating the dedicated buffer allocation to queues, the software rounds off any fractional dedicated buffer value to the closest lower full integer and programs this value in the hardware to avoid over allocation.

After allocating dedicated buffers to all configured queues, all QFX5000 platforms allocate any unused port dedicated buffers space to the first configured queue.

### Configuring Dedicated Port Buffer Allocation to Queues

In a port configuration that includes multiple forwarding class sets, with multiple forwarding classes mapped to multiple schedulers, the allocation of port dedicated buffers to queues depends on the mix of queues with buffer sizes configured as explicit percentages and queues configured with (or defaulted to) the *remainder* option.

The best way to demonstrate how using the percentage and remainder options affects dedicated port buffer allocation to queues is by showing an example of queue buffer allocation, and then showing how the queue buffer allocation changes when you add another forwarding class (queue) to the port.



[Table 107 on page 686](#) shows an initial configuration that includes four forwarding class sets, the five default forwarding classes (mapped to the five default queues for those forwarding classes), the buffer-size option configuration, and the resulting buffer allocation for each queue. [Table 108 on page 687](#) shows the same configuration after we add another forwarding class (best-effort-2, mapped to queue 1) to the best-effort forwarding class set. Comparing the buffer allocations in each table shows you how adding another queue affects buffer allocation when you use remainders and explicit percentages to configure the buffer allocation for different queues.

**Table 107: Egress Queue Dedicated Buffer Allocation (Example 1)**

| Forwarding Class Set (Priority Group) | Forwarding Class | Queue | Scheduler Buffer Size Configuration | Buffer Allocation per Queue (Percentage) |
|---------------------------------------|------------------|-------|-------------------------------------|------------------------------------------|
| fc-set-be                             | best-effort      | 0     | 10%                                 | 10%                                      |
| fc-set-lossless                       | fcoe             | 3     | 20%                                 | 20%                                      |
|                                       | no-loss          | 4     | 40%                                 | 40%                                      |
| fc-set-strict-high                    | network-control  | 7     | remainder                           | 15%                                      |
| fc-set-mcast                          | mcast            | 8     | remainder                           | 15%                                      |

In this first example, 70 percent of the egress port dedicated buffer pool is explicitly allocated to the best-effort, fcoe, and no-loss queues. The remaining 30 percent of the port dedicated buffer pool is split between the two queues that use the *remainder* option (network-control and mcast), so each queue receives 15 percent of the dedicated buffer pool.

Now we add another forwarding class (queue) to the best-effort priority group (fc-set-be) and configure it with a buffer size of *remainder* instead of configuring a specific percentage. Because a third queue now shares the remaining dedicated buffers, the queues that share the remainder receive fewer dedicated buffers, as shown in [Table 108 on page 687](#). The queues with explicitly configured percentages receive the configured percentage of dedicated buffers.

**Table 108: Egress Queue Dedicated Buffer Allocation with Another Remainder Queue (Example 2)**

| Priority Group (fc-set) | Forwarding Class | Queue | Scheduler Buffer Size Configuration | Buffer Allocation per Queue (Percentage) |
|-------------------------|------------------|-------|-------------------------------------|------------------------------------------|
| fc-set-be               | best-effort      | 0     | 10%                                 | 10%                                      |
|                         | best-effort-2    | 1     | remainder                           | 10%                                      |
| fc-set-lossless         | fcoe             | 3     | 20%                                 | 20%                                      |
|                         | no-loss          | 4     | 40%                                 | 40%                                      |
| fc-set-strict-high      | network-control  | 7     | remainder                           | 10%                                      |
| fc-set-mcast            | mcast            | 8     | remainder                           | 10%                                      |

The two tables show how the port divides the dedicated buffer space that remains after servicing the queues that have an explicitly configured percentage of dedicated buffer space.

### Trade-off Between Shared Buffer Space and Dedicated Buffer Space

The trade-off between shared buffer space and dedicated buffer space is:

- Shared buffers provide better absorption of traffic bursts because there is a larger pool of dynamic buffers that ports can use as needed to handle the bursts. However, all flows that exhaust their dedicated buffer space compete for the shared buffer pool. A larger shared buffer pool means a smaller dedicated buffer pool, and therefore more competition for the shared buffer pool because more flows exhaust their dedicated buffer allocation. Too much shared buffer space results in no single flow receiving very much shared buffer space, to maintain fairness when many flows contend for that space.
- Dedicated buffers provide guaranteed buffer space to each port. The larger the dedicated buffer pool, the less likely that congestion on one port affects traffic on another port, because the traffic does not need to use as much shared buffer space. However, less shared buffer space means less ability to dynamically absorb traffic bursts.

For optimal burst absorption, the switch needs enough dedicated buffer space to avoid persistent competition for the shared buffer space. When fewer flows compete for the shared buffers, the flows

that need shared buffer space to absorb bursts receive more of the shared buffer because fewer flows exhaust their dedicated buffer space.

The default configuration and the configurations recommended for different traffic scenarios allocate 100 percent of the user-configurable memory space to the global shared buffer pool because the amount of space reserved for dedicated buffers provides enough space to avoid persistent competition for dynamic shared buffers. This results in fewer flows competing for the shared buffers, so the competing flows receive more of the buffer space.

## Order of Buffer Consumption

The total buffer pool is divided into ingress and egress shared buffer pools and dedicated buffer pools. When traffic flows through the switch, the buffer space is used in a particular order that depends on the type of traffic.

On ingress, the order of buffer consumption is:

- Best-effort unicast traffic:
  1. Dedicated buffers
  2. Shared buffers
  3. Global headroom buffers (very small)
- Lossless unicast traffic:
  1. Dedicated buffers
  2. Shared buffers
  3. Lossless headroom buffers
  4. Global headroom buffers (very small)
- Multidestination traffic:
  1. Dedicated buffers
  2. Shared buffers
  3. Global headroom buffers (very small)

On egress, the order of buffer consumption is the same for unicast best-effort, lossless unicast, and multidestination traffic:

- Dedicated buffers
- Shared buffers

In all cases on all ports, the switch uses the dedicated buffer pool first and the shared buffer pool only after the dedicated buffer pool for the port or queue is exhausted. This reserves the maximum amount of dynamic shared buffer space to absorb traffic bursts.

## Default Buffer Pool Values

You can view the default or configured ingress and egress buffer pool values in KB units using the `show class-of-service shared-buffer operational` command. You can view the configured shared buffer pool values in percent units using the `show configuration class-of-service shared-buffer operational` command.

This section provides the default total buffer, shared buffer, and dedicated buffer values.

### Total Buffer Pool Size

The total buffer pool is common memory that has separate ingress and egress accounting, so the full buffer pool is available from both the ingress and egress perspective. The total buffer pool consists of the dedicated buffer space and the shared buffer space. The size of the total buffer pool is not user-configurable, but the allocation of buffer space to the dedicated and shared buffer pools is user-configurable.

On QFX3500 and QFX3600 switches, the combined total size of the ingress and egress buffer pools is approximately 9 MB (exactly 9360 KB).

On QFX5100, EX4600, and OCX Series switches, the combined total size of the ingress and egress buffer pools is approximately 12 MB (exactly 12480 KB).

On QFX5110 and QFX5200-32C switches, the combined total size of the ingress and egress buffer pools is approximately 16 MB.

On QFX5200-48Y switches, the combined total size of the ingress and egress buffer pools is approximately 22 MB.

On QFX5210 switches, the combined total size of the ingress and egress buffer pools is approximately 42 MB.

On QFX5220 switches, the combined total size of the ingress and egress buffer pools is approximately 64 MB.

### Shared Buffer Pool Default Values

Some switches have a larger shared buffer pool than other switches. However, the allocation of shared buffer space to the individual ingress and egress buffer pools is the same on a percentage basis, even though the absolute values are different. For example, the default ingress lossless buffer is 9 percent of the total shared ingress buffer space on all of the switches, even though the default absolute value of the ingress lossless buffer differs from switch to switch.

## Shared Ingress Buffer Default Values

Table 109 on page 690 shows the default ingress shared buffer allocation values in KB units for QFX5210 switches.

**Table 109: QFX5210 Switch Default Shared Ingress Buffer Values (KB)**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 29224                       | 2630.16         | 13150.80                 | 13443.04     |

Table 110 on page 690 shows the default ingress shared buffer allocation values in KB units for QFX5200-48Y switches.

**Table 110: QFX5200-48Y Switch Default Shared Ingress Buffer Values (KB)**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 19154.69                    | 1723.92         | 8619.61                  | 8811.16      |

Table 111 on page 690 shows the default ingress shared buffer allocation values in KB units for QFX5110 and QFX5200-32C switches.

**Table 111: QFX5110 and QFX5200-32C Switch Default Shared Ingress Buffer Values (KB)**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 11779.62                    | 1060.17         | 5300.83                  | 5418.63      |

Table 112 on page 690 shows the default ingress shared buffer allocation values in KB units for QFX5100, EX4600, and OCX Series switches.

**Table 112: QFX5100, EX4600, and OCX Series Switch Default Shared Ingress Buffer Values (KB)**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 9567.19 KB                  | 861.05 KB       | 4305.23 KB               | 4400.91 KB   |

[Table 113 on page 691](#) shows the default ingress shared buffer allocation values in KB units for QFX3500 and QFX3600 switches.

**Table 113: QFX3500 and QFX3600 Switch Default Shared Ingress Buffer Values (KB)**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 7202 KB                     | 648.18 KB       | 3240.9 KB                | 3312.92 KB   |

[Table 114 on page 691](#) shows the default ingress shared buffer allocation values as percentages for all switches. (If you change the default shared buffer allocation, you configure the change as a percentage.)

**Table 114: Default Shared Ingress Buffer Values (Percentage)**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 100%                        | 9%              | 45%                      | 46%          |

#### Shared Egress Buffer Default Values

[Table 115 on page 691](#) shows the default egress shared buffer allocation values in KB units for QFX5210 switches.

**Table 115: QFX5210 Switch Default Shared Egress Buffer Values (KB)**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 28080                      | 14040           | 8704.80      | 5335.20          |

[Table 116 on page 691](#) shows the default egress shared buffer allocation values in KB units for QFX5200-48Y switches.

**Table 116: QFX5200-48Y Switch Default Shared Egress Buffer Values (KB)**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 19115.69                   | 9557.84         | 5925.86      | 3631.98          |

[Table 117 on page 692](#) shows the default egress shared buffer allocation values in KB units for QFX5110 and QFX5200-32C switches.

**Table 117: QFX5110 and QFX5200-32C Switch Default Shared Egress Buffer Values (KB)**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 11232                      | 5616            | 3481.92      | 2134             |

**NOTE:** QFX5200-32C does not replicate all multicast streams when two or more downstream interface packet sizes are higher than ~6k and have an 1000pps packet ingress rate. This is because the number of working flows on QFX5200-32C is indirectly proportional to the packet size and directly proportional to available multicast shared buffers.

[Table 118 on page 692](#) shows the default egress shared buffer allocation values in KB units for QFX5100, EX4600, and OCX Series switches.

**Table 118: QFX5100, EX4600, and OCX Series Switch Default Shared Egress Buffer Values (KB)**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 8736 KB                    | 4368 KB         | 2708.16 KB   | 1659.84 KB       |

[Table 119 on page 692](#) shows the default egress shared buffer allocation values in KB units.

**Table 119: QFX3500 and QFX3600 Switch Default Shared Egress Buffer Values (KB)**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 6656 KB                    | 3328 KB         | 2063.36 KB   | 1264.64 KB       |

[Table 120 on page 693](#) shows the default egress shared buffer allocation values for all switches as percentages.

**Table 120: Default Shared Egress Buffer Values (Percentage)**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 100%                       | 50%             | 31%          | 19%              |

### Dedicated Buffer Pool Default Values

The system reserves ingress and egress dedicated buffer pools that are divided equally among the switch ports. By default, the system allocates 100 percent of the available unreserved buffer space to the shared buffer pool. If you reduce the percentage of available unreserved buffer space allocated to the shared buffer pool, the remaining unreserved buffer space is added to the dedicated buffer pool allocation. You configure the amount of dedicated buffer pool space by reducing (or increasing) the percentage of buffer space allocated to the shared buffer pool. You do not directly configure the dedicated buffer pool allocation.

[Table 121 on page 693](#) shows the default ingress and egress dedicated buffer pool values in KB units for QFX5210, QFX5200, QFX5110, QFX5100, QFX3500, QFX3600, EX4600, and OCX Series switches.

**Table 121: Default Ingress and Egress Dedicated Buffer Pool Values (KB) per Switch (**

| Dedicated Buffer Type | QFX5210 | QFX5200-48Y | QFX5110,<br>QFX5200-32C | QFX5100,<br>EX4600, OCX<br>Series | QFX3500,<br>QFX3600 |
|-----------------------|---------|-------------|-------------------------|-----------------------------------|---------------------|
| Ingress               | 14040   | 3373.50     | 4860.38                 | 2912.81                           | 2158                |
| Egress                | 15184   | 3412.50     | 5408                    | 3744                              | 2704                |

### Shared Buffer Configuration Recommendations for Different Network Traffic Scenarios

The way you configure the shared buffer pool depends on the mix of traffic on your network. This section provides shared buffer configuration recommendations for five basic network traffic scenarios:

- **Balanced traffic**—The network carries a balanced mix of unicast best-effort, lossless, and multicast traffic. (This is the default configuration.)
- **Best-effort unicast traffic**—The network carries mostly unicast best-effort traffic.



- Best-effort traffic with Ethernet PAUSE (IEEE 802.3X) enabled—The network carries mostly best-effort traffic with Ethernet PAUSE enabled on the links.
- Best-effort multicast traffic—The network carries mostly multicast best-effort traffic.
- Lossless traffic—The network carries mostly lossless traffic (traffic on which PFC is enabled).

**NOTE:** Lossless traffic is defined as traffic on which you enable PFC to ensure lossless transport. Lossless traffic does not refer to best-effort traffic on a link on which you enable Ethernet PAUSE. Start with the recommended profiles for each network traffic scenario, and adjust them if necessary for your network traffic conditions.

OCX Series switches do not support lossless transport or PFC. In this topic, references to lossless transport do not apply to OCX Series switches. OCX Series switches support symmetric Ethernet PAUSE.



**CAUTION:** Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete. This includes changing the default configuration to one of the recommended configurations.

Because you configure buffer allocations in percentages, the recommended allocations for each network traffic scenario are valid for all QFX Series switches, EX4600 switches, and OCX Series switches. Use one of the following recommended shared buffer configurations for your network traffic conditions. Start with a recommended configuration, then make small adjustments to the buffer allocations to fine-tune the buffers if necessary as described in ["Optimizing Buffer Configuration" on page 698](#).

### Balanced Traffic (Default Configuration)

The default shared buffer configuration is optimized for networks that carry a balanced mix of best-effort unicast, lossless, and multidestination (multicast, broadcast, and destination lookup fail) traffic. The default class-of-service (CoS) configuration is also optimized for networks that carry a balanced mix of traffic.

**NOTE:** On OCX Series switches, the default CoS configuration optimization does not include lossless traffic because OCX Series switches do not support lossless transport.

Except on OCX Series switches, we recommend that you use the default shared buffer configuration for networks that carry a balanced mix of traffic, especially if you are using the default CoS settings. [Table 122 on page 695](#) shows the default ingress shared buffer allocations:

**Table 122: Default Ingress Shared Buffer Configuration**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 100%                        | 9%              | 45%                      | 46%          |

[Table 123 on page 695](#) shows the default egress shared buffer allocations:

**Table 123: Default Egress Shared Buffer Configuration**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 100%                       | 50%             | 31%          | 19%              |

### Best-Effort Unicast Traffic

If your network carries mostly best-effort (lossy) unicast traffic, then the default shared buffer configuration allocates too much buffer space to support lossless transport. Instead of wasting those buffers, we recommend that you use the following ingress shared buffer settings (see [Table 124 on page 695](#)) and egress shared buffer settings (see [Table 125 on page 695](#)):

**Table 124: Recommended Ingress Shared Buffer Configuration for Networks with Mostly Best-Effort Unicast Traffic**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 100%                        | 5%              | 0%                       | 95%          |

**Table 125: Recommended Egress Shared Buffer Configuration for Networks with Mostly Best-Effort Unicast Traffic**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 100%                       | 5%              | 75%          | 20%              |

See [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic](#) for an example that shows you how to configure the recommended buffer settings shown in [Table 124 on page 695](#) and [Table 125 on page 695](#).

### Ethernet PAUSE Traffic

If your network carries mostly best-effort (lossy) traffic *and* enables Ethernet PAUSE on links, then the default shared buffer configuration allocates too much buffer space to the shared ingress buffer (Ethernet PAUSE traffic uses the dedicated buffers instead of shared buffers) and not enough space to the lossless-headroom buffers. We recommend that you use the following ingress shared buffer settings (see [Table 126 on page 696](#)) and egress shared buffer settings (see [Table 127 on page 696](#)):

**Table 126: Recommended Ingress Shared Buffer Configuration for Networks with Mostly Best-Effort Traffic and Ethernet PAUSE Enabled**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 70%                         | 5%              | 80%                      | 15%          |

**Table 127: Recommended Egress Shared Buffer Configuration for Networks with Mostly Best-Effort Traffic and Ethernet PAUSE Enabled**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 100%                       | 5%              | 75%          | 20%              |

See [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled](#) for an example that shows you how to configure the recommended buffer settings shown in [Table 124 on page 695](#) and [Table 125 on page 695](#).

### Best-Effort Multicast (Multidestination) Traffic

If your network carries mostly best-effort (lossy) multicast traffic, then the default shared buffer configuration allocates too much buffer space to support lossless transport. Instead of wasting those buffers, we recommend that you use the following ingress shared buffer settings (see [Table 128 on page 697](#)) and egress shared buffer settings (see [Table 129 on page 697](#)):

**Table 128: Recommended Ingress Shared Buffer Configuration for Networks with Mostly Best -Effort Multicast Traffic**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 100%                        | 5%              | 0%                       | 95%          |

**Table 129: Recommended Egress Shared Buffer Configuration for Networks with Mostly Best-Effort Multicast Traffic**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 100%                       | 5%              | 20%          | 75%              |

See [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic](#) for an example that shows you how to configure the recommended buffer settings shown in [Table 128 on page 697](#) and [Table 129 on page 697](#).

## Lossless Traffic

If your network carries mostly lossless traffic, then the default shared buffer configuration allocates too much buffer space to support best-effort traffic. Instead of wasting those buffers, we recommend that you use the following ingress shared buffer settings (see [Table 130 on page 697](#)) and egress shared buffer settings (see [Table 131 on page 698](#)):

**Table 130: Recommended Ingress Shared Buffer Configuration for Networks with Mostly Lossless Traffic**

| Total Shared Ingress Buffer | Lossless Buffer | Lossless-Headroom Buffer | Lossy Buffer |
|-----------------------------|-----------------|--------------------------|--------------|
| 100%                        | 15%             | 80%                      | 5%           |

**Table 131: Recommended Egress Shared Buffer Configuration for Networks with Mostly Lossless Traffic**

| Total Shared Egress Buffer | Lossless Buffer | Lossy Buffer | Multicast Buffer |
|----------------------------|-----------------|--------------|------------------|
| 100%                       | 90%             | 5%           | 5%               |

See [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic](#) for an example that shows you how to configure the recommended buffer settings shown in [Table 130 on page 697](#) and [Table 131 on page 698](#).

## Optimizing Buffer Configuration

Starting from the default configuration or from a recommended buffer configuration, you can further optimize the buffer allocation to best support the mix of traffic on your network. Adjust the settings gradually to fine-tune the shared buffer allocation. Use caution when adjusting the shared buffer configuration, not just when you fine-tune the ingress and egress buffer partitions, but also when you fine-tune the total ingress and egress shared buffer percentage. (Remember that if you allocate less than 100 percent of the available buffers to the shared buffers, the remaining buffers are added to the dedicated buffers). Tuning the buffers incorrectly can cause problems such as ingress port congestion.



**CAUTION:** Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete.

The relationship between the sizes of the ingress buffer pool and the egress buffer pool affects when and where packets are dropped. The buffer pool sizes include the shared buffers and the dedicated buffers. In general, if there are more ingress buffers than egress buffers, the switch can experience ingress port congestion because egress queues fill before ingress queues can empty.

Use the `show class-of-service shared-buffer` operational command to see the sizes in kilobytes (KB) of the dedicated and shared buffers and of the shared buffer partitions.

For best-effort traffic (unicast and multidestination), the combined ingress lossy shared buffer partition and ingress dedicated buffers must be *less than* the combined egress lossy and multicast shared buffer partitions plus the egress dedicated buffers. This prevents ingress port congestion by ensuring that egress best-effort buffers are deeper than ingress best-effort buffers, and ensures that if packets are dropped, they are dropped at the egress queues. (Packets dropping at the ingress prevents the egress schedulers from working properly.)

For lossless traffic (traffic on which you enable PFC), the combined ingress lossless shared buffer partition and a reasonable portion of the ingress headroom buffer partition, plus the dedicated buffers,

must be *less than* the total egress lossless shared buffer partition and dedicated buffers. (A reasonable portion of the ingress headroom buffer is approximately 20 to 25 percent of the buffer space, but this varies depending on how much buffer headroom is required to support the lossless traffic.) When these conditions are met, if there is ingress port congestion, the ingress port congestion triggers PFC on the ingress port to prevent packet loss. If the total lossless ingress buffers exceed the total lossless egress buffers, packets could be dropped at the egress instead of PFC being applied at the ingress to prevent packet loss.

**NOTE:** If you commit a buffer configuration for which the switch does not have sufficient resources, the switch might log an error instead of returning a commit error. In that case, a syslog message is displayed on the console. For example:

```
user@host# commit
configuration check succeeds
```

```
Message from syslogd@host at Jun 13 11:11:10 ...
host dc-pfe: Not enough Ingress Lossless headroom.(Already allocated more). Dedicated : 14340
Lossy : 47100 Lossless 4239 Headroom 21195 Avail : 20781
commit complete
```

If the buffer configuration commits but you receive a syslog message that indicates the configuration cannot be implemented, you can:

- Reconfigure the buffers or reconfigure other parameters (for example, the PFC configuration, which affects the need for lossless headroom buffers and lossless buffers—the more priorities you pause, the more lossless and lossless headroom buffer space you need), then attempt the commit operation again.
- Roll back the switch to the last successful configuration.

If you receive a syslog message that says the buffer configuration cannot be implemented, you must take corrective action. If you do not fix the configuration or roll back to a previous successful configuration, the system behavior is unpredictable.

## General Buffer Configuration Rules and Considerations

Keep the following rules and considerations in mind when you configure the buffers:

- Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete.
- If you configure the ingress or egress shared buffer percentages as less than 100 percent, the remaining percentage of buffer space is added to the dedicated buffer pool.

- The sum of all of the ingress shared buffer partitions must equal 100 percent. Each partition must be configured with a value of at least 5 percent except the lossless headroom buffer, which can have a value of 0 percent.
- The sum of all of the egress shared buffer partitions must equal 100 percent. Each partition must be configured with a value of at least 5 percent.
- Lossless and lossless headroom shared buffers serve traffic on which you enable PFC, and do not serve traffic subject to Ethernet PAUSE.
- The switch uses the dedicated buffer pool first and the shared buffer pool only after the dedicated buffer pool for a port or queue is exhausted.
- Too little dedicated buffer space results in too much competition for shared buffer space.
- Too much dedicated buffer space results in poorer burst absorption because there is less available shared buffer space.
- Always check the syslog messages after you commit a new buffer configuration.
- The optimal buffer configuration for your network depends on the types of traffic on the network. If your network carries less traffic of a certain type (for example, lossless traffic), then you can reduce the size of the buffers allocated to that type of traffic (for example, you can reduce the sizes of the lossless and lossless headroom buffers).

## RELATED DOCUMENTATION

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic*

*Example: Configuring Queue Schedulers*

*Configuring Global Ingress and Egress Shared Buffers*

## Configuring Global Ingress and Egress Shared Buffers

Although the switch reserves some buffer space to ensure a minimum memory allocation for ports and queues, you can configure how the system uses the rest of the buffer space to optimize the buffer allocation for your particular mix of network traffic. The global shared buffer pool is memory space that all of the ports on the switch share dynamically as they need buffers. You can allocate global shared memory space to different types of ingress and egress buffers to better support different mixes of network traffic.



**CAUTION:** Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete.

Use the default shared buffer settings (for a network with a balanced mix of lossless, best-effort, and multicast traffic) or one of the recommended shared buffer configurations for your mix of network traffic (mostly best-effort unicast traffic, mostly best-effort traffic on links enabled for Ethernet PAUSE, mostly multicast traffic, or mostly lossless traffic). Either the default configuration or one of the recommended configurations provides a buffer allocation that satisfies the needs of most networks.

After starting from one of the recommended configurations, you can fine-tune the shared buffer settings, but do so with caution to prevent traffic loss due to buffer misconfiguration.

You can configure the percentage of available (user-configurable) buffer space allocated to the global shared buffers. Any space that you do not allocate to the global shared buffer pool is added to the dedicated buffer pool. The default configuration allocates 100 percent of the available buffer space to the global shared buffers.

You can partition the ingress and egress shared buffer pools to allocate more buffers to the types of traffic your network predominantly carries, and fewer buffers to other traffic. From the buffer space allocated to the ingress shared buffer pool, you can allocate space to:

- Lossless buffers—Percentage of shared buffer pool for all lossless ingress traffic. The minimum value for the lossless buffers is 5 percent.
- Lossless headroom buffers—Percentage of shared buffer pool for packets received while a pause is asserted. If Ethernet PAUSE is configured on a port or if priority-based flow control (PFC) is configured on priorities on a port, when the port sends a pause message to the connected peer, the port uses the headroom buffers to store the packets that arrive between the time the port sends the pause message and the time the last packet arrives after the peer pauses traffic. The minimum value for the lossless headroom buffers is 0 (zero) percent. (Lossless headroom buffers are the only buffers that can have a minimum value of less than 5 percent.)



**NOTE:** On a QFX Virtual Chassis and an EX4600/EX4650 Virtual Chassis, the minimum value for the lossless headroom buffer is 3 percent.

- Lossy buffers—Percentage of shared buffer pool for all best-effort ingress traffic (best-effort unicast, multideestination, and strict-high priority traffic). The minimum value for the lossy buffers is 5 percent.

The combined percentage values of the ingress lossless, lossless headroom, and lossy buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All ingress buffer partitions must be explicitly configured, even when the lossless headroom buffer partition has a value of 0 (zero) percent.

From the buffer space allocated to the egress shared buffer pool, you can allocate space to:

- Lossless buffers—Percentage of shared buffer pool for all lossless egress queues. The minimum value for the lossless buffers is 5 percent.
- Lossy buffers—Percentage of shared buffer pool for all best-effort egress queues (best-effort unicast and strict-high priority queues). The minimum value for the lossy buffers is 5 percent.
- Multicast buffers—Percentage of shared buffer pool for all multideestination (multicast, broadcast, and destination lookup fail) egress queues. The minimum value for the multicast buffers is 5 percent.

The combined percentage values of the egress lossless, lossy, and multicast buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All egress buffer partitions must be explicitly configured and must have a value of at least 5 percent.

To configure the shared buffer allocation and partitioning using the CLI:

1. Configure the percentage of available (nonreserved) buffers used for the ingress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set ingress percent percent
```

2. Configure the global ingress buffer partitions for lossless, lossless-headroom, and lossy traffic:

```
[edit class-of-service shared-buffer]
user@switch# set ingress buffer-partition lossless percent percent
user@switch# set ingress buffer-partition lossless-headroom percent percent
user@switch# set ingress buffer-partition lossy percent percent
```

3. Configure the percentage of available (nonreserved) buffers used for the egress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set egress percent percent
```

4. Configure the global egress buffer partitions for lossless, lossy, and multicast queues:

```
[edit class-of-service shared-buffer]
user@switch# set egress buffer-partition lossless percent percent
user@switch# set egress buffer-partition lossy percent percent
user@switch# set egress buffer-partition multicast percent percent
```

## RELATED DOCUMENTATION

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic*

*Understanding CoS Buffer Configuration*

## Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic

### IN THIS SECTION

- [Requirements | 704](#)
- [Overview | 704](#)
- [Configuration | 706](#)
- [Verification | 709](#)

Although the switch reserves some buffer space to ensure a minimum memory allocation for ports and queues, you can configure how the system uses the rest of the buffer space to optimize the buffer allocation for your particular mix of network traffic.

This example shows you the recommended configuration of the global shared buffer pool to support a network that carries mostly best-effort (lossy) unicast traffic. The global shared buffer pool is memory space that all of the ports on the switch share dynamically as they need buffers. You can allocate global shared memory space to different types of buffers to better support different mixes of network traffic.



**CAUTION:** Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete.

Use the default shared buffer settings (for a network with a balanced mix of lossless, best effort, and multicast traffic) or one of the recommended shared buffer configurations for your mix of network traffic (mostly best-effort unicast traffic, mostly best-effort traffic on links enabled for Ethernet PAUSE, mostly multicast traffic, or mostly lossless traffic). Either the default configuration or one of the recommended configurations provides a buffer allocation that satisfies the needs of most networks.

**NOTE:** OCX Series switches do not support lossless transport.

After starting from the recommended configuration, you can fine-tune the shared buffer settings, but do so with caution to prevent traffic loss due to buffer misconfiguration.

## Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 12.3 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

## Overview

### IN THIS SECTION

- [Topology | 706](#)

You can configure the percentage of available (user-configurable) buffer space allocated to the global shared buffers. Any space that you do not allocate to the global shared buffer pool is added to the dedicated buffer pool. The default configuration allocates 100 percent of the available buffer space to the global shared buffers.

You can partition the ingress and egress shared buffer pools to allocate more buffers to the types of traffic your network predominantly carries, and fewer buffers to other traffic. From the buffer space allocated to the ingress shared buffer pool, you can allocate space to:

- **Lossless buffers**—Percentage of shared buffer pool for all lossless ingress traffic. The minimum value for the lossless buffers is 5 percent.
- **Lossless headroom buffers**—Percentage of shared buffer pool for packets received while a pause is asserted. If Ethernet PAUSE is configured on a port or if priority-based flow control (PFC) is configured on priorities on a port, when the port sends a pause message to the connected peer, the port uses the headroom buffers to store the packets that arrive between the time the port sends the pause message and the time the last packet arrives after the peer pauses traffic. The minimum value for the lossless headroom buffers is 0 (zero) percent. (Lossless headroom buffers are the only buffers that can have a minimum value of less than 5 percent.)
- **Lossy buffers**—Percentage of shared buffer pool for all best-effort ingress traffic (best-effort unicast, multdestination, and strict-high priority traffic). The minimum value for the lossy buffers is 5 percent.

The combined percentage values of the ingress lossless, lossless headroom, and lossy buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All ingress buffer partitions must be explicitly configured, even when the lossless headroom buffer partition has a value of 0 (zero) percent.

From the buffer space allocated to the egress shared buffer pool, you can allocate space to:

- **Lossless buffers**—Percentage of shared buffer pool for all lossless egress queues. The minimum value for the lossless buffers is 5 percent.
- **Lossy buffers**—Percentage of shared buffer pool for all best-effort egress queues (best-effort unicast, and strict-high priority queues). The minimum value for the lossy buffers is 5 percent.
- **Multicast buffers**—Percentage of shared buffer pool for all multdestination (multicast, broadcast, and destination lookup fail) egress queues. The minimum value for the multicast buffers is 5 percent.

The combined percentage values of the egress lossless, lossy, and multicast buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All egress buffer partitions must be explicitly configured and must have a value of at least 5 percent.

To configure the shared buffers to support a network that carries mostly best-effort unicast traffic, more buffer space needs to be allocated to lossy buffers, and less buffer space should be allocated to lossless

buffers. This example shows you how to configure the global shared buffer pool allocation that we recommend to support a network that carries mostly unicast traffic.

Topology

Table 132 on page 706 shows the configuration components for this example.

**Table 132: Components of the Recommended Shared Buffer Configuration for Best-Effort Unicast Network Topologies**

| Component             | Settings                                                                                                                                                                                                                                                                                                                                                                                                                             |
|-----------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Hardware              | QFX3500 switch                                                                                                                                                                                                                                                                                                                                                                                                                       |
| Ingress shared buffer | Percentage of available ingress buffer space allocated to the ingress shared buffer: 100%<br><br>Percentage of ingress buffer space allocated to lossless traffic (lossless buffer partition): 5%<br><br>Percentage of ingress buffer space allocated to lossless headroom traffic (lossless-headroom buffer partition): 0%<br><br>Percentage of ingress buffer space allocated to best-effort traffic (lossy buffer partition): 95% |
| Egress shared buffer  | Percentage of available egress buffer space allocated to the egress shared buffer: 100%<br><br>Percentage of egress buffer space allocated to lossless queues (lossless buffer partition): 5%<br><br>Percentage of egress buffer space allocated to best-effort queues (lossy buffer partition): 75%<br><br>Percentage of egress buffer space allocated to multicast traffic (multicast buffer partition): 20%                       |

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 707](#)
- [Configuring the Global Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic | 707](#)
- [Results | 708](#)

## CLI Quick Configuration

To quickly configure the recommended shared buffer settings for networks that carry mostly best-effort unicast traffic, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the `[edit class-of-service shared-buffer]` hierarchy level:

```
[edit class-of-service shared-buffer]
set ingress percent 100
set ingress buffer-partition lossless percent 5
set ingress buffer-partition lossless-headroom percent 0
set ingress buffer-partition lossy percent 95
set egress percent 100
set egress buffer-partition lossless percent 5
set egress buffer-partition lossy percent 75
set egress buffer-partition multicast percent 20
```

## Configuring the Global Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic

### Step-by-Step Procedure

To configure the global ingress and egress shared buffer allocations and partitions for a network that carries mostly best-effort unicast traffic:

1. Configure the percentage of available (nonreserved) buffers used for the ingress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set ingress percent 100
```

2. Configure the global ingress buffer partitions for lossless, lossless-headroom, and lossy traffic:

```
[edit class-of-service shared-buffer]
user@switch# set ingress buffer-partition lossless percent 5
user@switch# set ingress buffer-partition lossless-headroom percent 0
user@switch# set ingress buffer-partition lossy percent 95
```

3. Configure the percentage of available (nonreserved) buffers used for the egress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set egress percent 100
```

4. Configure the global egress buffer partitions for lossless, lossy, and multicast queues:

```
[edit class-of-service shared-buffer]
user@switch# set egress buffer-partition lossless percent 5
user@switch# set egress buffer-partition lossy percent 75
user@switch# set egress buffer-partition multicast percent 20
```

## Results

Display the results of the configuration:

```
root@dcbg-tp-pa-02> show configuration class-of-service shared-buffer
ingress {
 percent 100;
 buffer-partition lossless {
 percent 5;
 }
 buffer-partition lossy {
 percent 95;
 }
 buffer-partition lossless-headroom {
 percent 0;
 }
}
egress {
 percent 100;
 buffer-partition lossless {
 percent 5;
 }
 buffer-partition lossy {
 percent 75;
 }
 buffer-partition multicast {
```

```
 percent 20;
 }
}
```

Verification

IN THIS SECTION

[Verifying the Shared Buffer Configuration | 709](#)

Verify that you correctly configured the shared buffer.

Verifying the Shared Buffer Configuration

Purpose

Verify that the ingress and egress global shared buffer pools are correctly configured and partitioned among the shared buffer types.

Action

List the global shared buffer configuration using the operational mode command `show class-of-service shared-buffer`:

```
user@switch> show class-of-service shared-buffer
root@dcbg-tp-pa-02> show class-of-service shared-buffer
Ingress:
 Total Buffer : 9360.00 KB
 Dedicated Buffer : 2158.00 KB
 Shared Buffer : 7202.00 KB
 Lossless : 360.10 KB
 Lossless Headroom : 0.00 KB
 Lossy : 6841.90 KB

Lossless Headroom Utilization:
Node Device Total Used Free
0 0.00 KB 0.00 KB 0.00 KB
```



**Egress:**

```

Total Buffer : 9360.00 KB
Dedicated Buffer : 2704.00 KB
Shared Buffer : 6656.00 KB
 Lossless : 332.80 KB
 Multicast : 1331.20 KB
 Lossy : 4992.00 KB

```

**Meaning**

The `show class-of-service shared-buffer operational` command shows all of the ingress and egress global shared buffer settings, including the buffer partitioning.

For the ingress shared buffers, the command output shows:

- The total switch buffer pool is 9360 KB (9 MB).
- The dedicated buffer pool is 2158 KB. This is the size of the global ingress dedicated buffer pool when you configure the ingress shared buffer pool as 100 percent of the available (user-configurable) buffer space. This is the minimum size of the reserved, ingress dedicated ingress buffer pool (not user-configurable). If you configure the shared buffer as less than 100 percent of the available buffer pool, the remaining buffer space is added to the dedicated buffer pool.
- With the ingress shared buffer pool configured as 100 percent of the available buffers, the total size of the ingress shared buffer pool is 7202 KB.
- The ingress shared buffer pool is partitioned to allocate:
  - 360.10 KB to lossless traffic
  - No space to lossless headroom traffic
  - 6841.90 KB to lossy unicast traffic
- The Lossless Headroom Utilization field shows how much of the buffer space reserved for paused traffic is used. Because the lossless headroom buffer partition is set to 0 (zero) percent, the total amount of lossless headroom buffer space is 0 KB; therefore the amount of used and free lossless headroom buffer space is also 0 KB.

For the egress shared buffers, the command output shows:

- The total switch buffer pool is 9360 KB (9 MB).
- The dedicated buffer pool is 2704 KB. This is the size of the global egress dedicated buffer pool when you configure the egress shared buffer pool as 100 percent of the available (user-configurable) buffer space. This is the minimum size of the reserved, egress dedicated buffer pool (not user-

configurable). If you configure the shared buffer as less than 100 percent of the available buffer pool, the remaining buffer space is added to the dedicated buffer pool.

- With the egress shared buffer pool configured as 100 percent of the available buffers, the total size of the egress shared buffer pool is 6656 KB. This is less than the ingress shared buffer pool because the switch reserves more egress dedicated buffer space than ingress dedicated buffer space. (More dedicated buffer space means less shared buffer space, and more shared buffer space means less dedicated buffer space.)
- The egress shared buffer pool is partitioned to allocate:
  - 332.80 KB to lossless traffic
  - 1331.20 KB to multicast traffic
  - 4992 KB to lossy unicast traffic

**NOTE:** The output values are valid for QFX3500 and QFX3600 switches. QFX5100, EX4600, and OCX Series switches have larger buffers (12 MB instead of 9 MB), so the total buffer size and the sizes of each buffer partition are larger on those switches.

## RELATED DOCUMENTATION

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic*

*Configuring Global Ingress and Egress Shared Buffers*

*Understanding CoS Buffer Configuration*

## Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled

### IN THIS SECTION

- [Requirements | 713](#)
- [Overview | 713](#)
- [Configuration | 715](#)
- [Verification | 718](#)

Although the switch reserves some buffer space to ensure a minimum memory allocation for ports and queues, you can configure how the system uses the rest of the buffer space to optimize the buffer allocation for your particular mix of network traffic.

This example shows you the recommended configuration of the global shared buffer pool to support a network that carries mostly best-effort (lossy) traffic on links with Ethernet PAUSE (IEEE 802.3X) enabled.

**NOTE:** OCX Series switches support symmetric Ethernet PAUSE flow control, but do not support asymmetric Ethernet PAUSE flow control.

The global shared buffer pool is memory space that all of the ports on the switch share dynamically as they need buffers. You can allocate global shared memory space to different types of buffers to better support different mixes of network traffic.



**CAUTION:** Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete.

Use the default shared buffer settings (for a network with a balanced mix of lossless, best effort, and multicast traffic) or one of the recommended shared buffer configurations for your mix of network traffic (mostly best-effort unicast traffic, mostly best-effort traffic on links enabled for Ethernet PAUSE, mostly multicast traffic, or mostly lossless traffic). Either the default configuration or one of the recommended configurations provides a buffer allocation that satisfies the needs of most networks.

After starting from the recommended configuration, you can fine-tune the shared buffer settings, but do so with caution to prevent traffic loss due to buffer misconfiguration.

## Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 12.3 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

## Overview

### IN THIS SECTION

- [Topology | 714](#)

You can configure the percentage of available (user-configurable) buffer space allocated to the global shared buffers. Any space that you do not allocate to the global shared buffer pool is added to the dedicated buffer pool. The default configuration allocates 100 percent of the available buffer space to the global shared buffers.

You can partition the ingress and egress shared buffer pools to allocate more buffers to the types of traffic your network predominantly carries, and fewer buffers to other traffic. From the buffer space allocated to the ingress shared buffer pool, you can allocate space to:

- Lossless buffers—Percentage of shared buffer pool for all lossless ingress traffic. The minimum value for the lossless buffers is 5 percent.
- Lossless headroom buffers—Percentage of shared buffer pool for packets received while a pause is asserted. If Ethernet PAUSE is configured on a port or if priority-based flow control (PFC) is configured on priorities on a port, when the port sends a pause message to the connected peer, the port uses the headroom buffers to store the packets that arrive between the time the port sends the pause message and the time the last packet arrives after the peer pauses traffic. The minimum value for the lossless headroom buffers is 0 (zero) percent. (Lossless headroom buffers are the only buffers that can have a minimum value of less than 5 percent.)

**NOTE:** OCX Series switches do not support PFC.

- Lossy buffers—Percentage of shared buffer pool for all best-effort ingress traffic (best-effort unicast, multidestination, and strict-high priority traffic). The minimum value for the lossy buffers is 5 percent.

The combined percentage values of the ingress lossless, lossless headroom, and lossy buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All ingress buffer partitions must be explicitly configured, even when the lossless headroom buffer partition has a value of 0 (zero) percent.

From the buffer space allocated to the egress shared buffer pool, you can allocate space to:

- Lossless buffers—Percentage of shared buffer pool for all lossless egress queues. The minimum value for the lossless buffers is 5 percent.
- Lossy buffers—Percentage of shared buffer pool for all best-effort egress queues (best-effort unicast and strict-high priority queues). The minimum value for the lossy buffers is 5 percent.
- Multicast buffers—Percentage of shared buffer pool for all multidestination (multicast, broadcast, and destination lookup fail) egress queues. The minimum value for the multicast buffers is 5 percent.

The combined percentage values of the egress lossless, lossy, and multicast buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All egress buffer partitions must be explicitly configured and must have a value of at least 5 percent.

To configure the shared buffers to support a network that carries mostly best-effort traffic on links enabled for Ethernet PAUSE, more buffer space needs to be allocated to ingress dedicated port buffers, and less buffer space should be allocated to ingress shared buffers. Also, more buffer space needs to be allocated to lossless-headroom buffers, and less space to ingress lossy buffers. This example shows you how to configure the global shared buffer pool allocation that we recommend to support a network that carries mostly best-effort traffic on links enabled for Ethernet PAUSE.

Topology

Table 133 on page 714 shows the configuration components for this example.

**Table 133: Components of the Recommended Shared Buffer Configuration for Best-Effort Network Topologies with Links Enabled for Ethernet PAUSE**

| Component | Settings       |
|-----------|----------------|
| Hardware  | QFX3500 switch |

**Table 133: Components of the Recommended Shared Buffer Configuration for Best-Effort Network Topologies with Links Enabled for Ethernet PAUSE *(Continued)***

| Component             | Settings                                                                                                                                                                                                                                                                                                                                                                                                                                    |
|-----------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Ingress shared buffer | <p>Percentage of available ingress buffer space allocated to the ingress shared buffer: 70%</p> <p>Percentage of ingress buffer space allocated to lossless traffic (lossless buffer partition): 5%</p> <p>Percentage of ingress buffer space allocated to lossless headroom traffic (lossless-headroom buffer partition): 80%</p> <p>Percentage of ingress buffer space allocated to best-effort traffic (lossy buffer partition): 15%</p> |
| Egress shared buffer  | <p>Percentage of available egress buffer space allocated to the egress shared buffer: 100%</p> <p>Percentage of egress buffer space allocated to lossless queues (lossless buffer partition): 5%</p> <p>Percentage of egress buffer space allocated to best-effort queues (lossy buffer partition): 75%</p> <p>Percentage of egress buffer space allocated to multicast traffic (multicast buffer partition): 20%</p>                       |

## Configuration

### IN THIS SECTION

- [CLI Quick Configuration | 715](#)
- [Configuring the Global Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links Enabled for Ethernet PAUSE | 716](#)
- [Results | 717](#)

### CLI Quick Configuration

To quickly configure the recommended shared buffer settings for networks that carry mostly best-effort unicast traffic, copy the following commands, paste them in a text file, remove line breaks, change

variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit class-of-service shared-buffer] hierarchy level:

```
[edit class-of-service shared-buffer]
set ingress percent 70
set ingress buffer-partition lossless percent 5
set ingress buffer-partition lossless-headroom percent 80
set ingress buffer-partition lossy percent 15
set egress percent 100
set egress buffer-partition lossless percent 5
set egress buffer-partition lossy percent 75
set egress buffer-partition multicast percent 20
```

## Configuring the Global Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links Enabled for Ethernet PAUSE

### Step-by-Step Procedure

To configure the global ingress and egress shared buffer allocations and partitions:

1. Configure the percentage of available (nonreserved) buffers used for the ingress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set ingress percent 70
```

2. Configure the global ingress buffer partitions for lossless, lossless-headroom, and lossy traffic:

```
[edit class-of-service shared-buffer]
user@switch# set ingress buffer-partition lossless percent 5
user@switch# set ingress buffer-partition lossless-headroom percent 80
user@switch# set ingress buffer-partition lossy percent 15
```

3. Configure the percentage of available (nonreserved) buffers used for the egress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set egress percent 100
```

4. Configure the global egress buffer partitions for lossless, lossy, and multicast queues:

```
[edit class-of-service shared-buffer]
user@switch# set egress buffer-partition lossless percent 5
user@switch# set egress buffer-partition lossy percent 75
user@switch# set egress buffer-partition multicast percent 20
```

## Results

Display the results of the configuration:

```
root@dcbg-tp-pa-02> show configuration class-of-service shared-buffer
ingress {
 percent 70;
 buffer-partition lossless {
 percent 5;
 }
 buffer-partition lossy {
 percent 15;
 }
 buffer-partition lossless-headroom {
 percent 80;
 }
}
egress {
 percent 100;
 buffer-partition lossless {
 percent 5;
 }
 buffer-partition lossy {
 percent 75;
 }
 buffer-partition multicast {
```



```

 percent 20;
 }
}

```

## Verification

### IN THIS SECTION

- [Verifying the Shared Buffer Configuration | 718](#)

Verify that you correctly configured the shared buffer.

### Verifying the Shared Buffer Configuration

#### Purpose

Verify that the ingress and egress global shared buffer pools are correctly configured and partitioned among the shared buffer types.

#### Action

List the global shared buffer configuration using the operational mode command `show class-of-service shared-buffer`:

```

user@switch> show class-of-service shared-buffer
root@dcbg-tp-pa-02> show class-of-service shared-buffer
Ingress:
 Total Buffer : 9360.00 KB
 Dedicated Buffer : 4318.60 KB
 Shared Buffer : 5041.40 KB
 Lossless : 252.07 KB
 Lossless Headroom : 4033.12 KB
 Lossy : 756.21 KB

Egress:
 Total Buffer : 9360.00 KB
 Dedicated Buffer : 2704.00 KB
 Shared Buffer : 6656.00 KB

```

|           |   |            |
|-----------|---|------------|
| Lossless  | : | 332.80 KB  |
| Multicast | : | 1331.20 KB |
| Lossy     | : | 4992.00 KB |

## Meaning

The `show class-of-service shared-buffer operational` command shows all of the ingress and egress global shared buffer settings, including the buffer partitioning.

For the ingress shared buffers, the command output shows:

- The total switch buffer pool is 9360 KB (9 MB).
- The dedicated buffer pool is 4318.6 KB. This is the size of the global ingress dedicated buffer pool when you configure the ingress shared buffer pool as 70 percent of the available (user-configurable) buffer space.
- With the ingress shared buffer pool configured as 70 percent of the available buffers, the total size of the ingress shared buffer pool is 5041.4 KB.
- The ingress shared buffer pool is partitioned to allocate:
  - 252.07 KB to lossless traffic
  - 4033.12 KB to lossless headroom traffic
  - 756.21 KB to lossy unicast traffic

For the egress shared buffers, the command output shows:

- The total switch buffer pool is 9360 KB (9 MB).
- The dedicated buffer pool is 2704 KB. This is the size of the global egress dedicated buffer pool when you configure the egress shared buffer pool as 100 percent of the available (user-configurable) buffer space. This is the minimum size of the reserved, egress dedicated buffer pool (not user-configurable). If you configure the shared buffer as less than 100 percent of the available buffer pool, the remaining buffer space is added to the dedicated buffer pool.
- With the egress shared buffer pool configured as 100 percent of the available buffers, the total size of the egress shared buffer pool is 6656 KB. This is less than the ingress shared buffer pool because the switch reserves more egress dedicated buffer space than ingress dedicated buffer space. (More dedicated buffer space means less shared buffer space, and more shared buffer space means less dedicated buffer space.)
- The egress shared buffer pool is partitioned to allocate:
  - 332.80 KB to lossless traffic

- 1331.20 KB to multicast traffic
- 4992 KB to lossy unicast traffic

**NOTE:** The output values are valid for QFX3500 and QFX3600 switches. QFX5100, EX4600, and OCX Series switches have larger buffers (12 MB instead of 9 MB), so the total buffer size and the sizes of each buffer partition are larger on those switches.

## RELATED DOCUMENTATION

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic*

*Configuring Global Ingress and Egress Shared Buffers*

*Understanding CoS Buffer Configuration*

## Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic

### IN THIS SECTION

- [Requirements | 721](#)
- [Overview | 721](#)
- [Configuration | 723](#)
- [Verification | 726](#)

Although the switch reserves some buffer space to ensure a minimum memory allocation for ports and queues, you can configure how the system uses the rest of the buffer space to optimize the buffer allocation for your particular mix of network traffic.

This example shows you the recommended configuration of the global shared buffer pool to support a network that carries mostly multicast traffic. The global shared buffer pool is memory space that all of the ports on the switch share dynamically as they need buffers. You can allocate global shared memory space to different types of buffers to better support different mixes of network traffic.



**CAUTION:** Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete.

Use the default shared buffer settings (for a network with a balanced mix of lossless, best effort, and multicast traffic) or one of the recommended shared buffer configurations for your mix of network traffic (mostly best-effort unicast traffic, mostly best-effort traffic on links enabled for Ethernet PAUSE, mostly multicast traffic, or mostly lossless traffic). Either the default configuration or one of the recommended configurations provides a buffer allocation that satisfies the needs of most networks.

After starting from the recommended configuration, you can fine-tune the shared buffer settings, but do so with caution to prevent traffic loss due to buffer misconfiguration.

## Requirements

This example uses the following hardware and software components:

- One switch (this example was tested on a Juniper Networks QFX3500 Switch)
- Junos OS Release 12.3 or later for the QFX Series or Junos OS Release 14.1X53-D20 or later for the OCX Series

## Overview

### IN THIS SECTION

- [Topology | 723](#)

You can configure the percentage of available (user-configurable) buffer space allocated to the global shared buffers. Any space that you do not allocate to the global shared buffer pool is added to the dedicated buffer pool. The default configuration allocates 100 percent of the available buffer space to the global shared buffers.

You can partition the ingress and egress shared buffer pools to allocate more buffers to the types of traffic your network predominantly carries, and fewer buffers to other traffic. From the buffer space allocated to the ingress shared buffer pool, you can allocate space to:

- **Lossless buffers**—Percentage of shared buffer pool for all lossless ingress traffic. The minimum value for the lossless buffers is 5 percent.
- **Lossless headroom buffers**—Percentage of shared buffer pool for packets received while a pause is asserted. If Ethernet PAUSE is configured on a port or if priority-based flow control (PFC) is configured on priorities on a port, when the port sends a pause message to the connected peer, the port uses the headroom buffers to store the packets that arrive between the time the port sends the pause message and the time the last packet arrives after the peer pauses traffic. The minimum value for the lossless headroom buffers is 0 (zero) percent. (Lossless headroom buffers are the only buffers that can have a minimum value of less than 5 percent.)
- **Lossy buffers**—Percentage of shared buffer pool for all best-effort ingress traffic (best-effort unicast, multdestination, and strict-high priority traffic). The minimum value for the lossy buffers is 5 percent.

**NOTE:** For virtual chassis deployments, you cannot configure virtual lossless headroom buffers with 0% value. You need a minimum buffer value of 5% for 2 VCP ports and if there are more ports, more buffers are required to configure lossless headroom partitions.

The combined percentage values of the ingress lossless, lossless headroom, and lossy buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All ingress buffer partitions must be explicitly configured, even when the lossless headroom buffer partition has a value of 0 (zero) percent.

From the buffer space allocated to the egress shared buffer pool, you can allocate space to:

- **Lossless buffers**—Percentage of shared buffer pool for all lossless egress queues. The minimum value for the lossless buffers is 5 percent.
- **Lossy buffers**—Percentage of shared buffer pool for all best-effort egress queues (best-effort unicast, and strict-high priority queues). The minimum value for the lossy buffers is 5 percent.
- **Multicast buffers**—Percentage of shared buffer pool for all multdestination (multicast, broadcast, and destination lookup fail) egress queues. The minimum value for the multicast buffers is 5 percent.

The combined percentage values of the egress lossless, lossy, and multicast buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All egress buffer partitions must be explicitly configured and must have a value of at least 5 percent.

To configure the shared buffers to support a network that carries mostly multicast traffic, more buffer space needs to be allocated to lossy buffers, less buffer space should be allocated to lossless buffers, and more space needs to be allocated to egress multicast buffers. This example shows you how to configure the global shared buffer pool allocation that we recommend to support a network that carries mostly multicast traffic.

# Topology

Table 134 on page 723 shows the configuration components for this example.

**Table 134: Components of the Recommended Shared Buffer Configuration for Multicast Network Topologies**

| Component             | Settings                                                                                                                                                                                                                                                                                                                                                                                                                                    |
|-----------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Hardware              | QFX3500 switch                                                                                                                                                                                                                                                                                                                                                                                                                              |
| Ingress shared buffer | <p>Percentage of available ingress buffer space allocated to the ingress shared buffer: 100%</p> <p>Percentage of ingress buffer space allocated to lossless traffic (lossless buffer partition): 5%</p> <p>Percentage of ingress buffer space allocated to lossless headroom traffic (lossless-headroom buffer partition): 0%</p> <p>Percentage of ingress buffer space allocated to best-effort traffic (lossy buffer partition): 95%</p> |
| Egress shared buffer  | <p>Percentage of available egress buffer space allocated to the egress shared buffer: 100%</p> <p>Percentage of egress buffer space allocated to lossless queues (lossless buffer partition): 5%</p> <p>Percentage of egress buffer space allocated to best-effort queues (lossy buffer partition): 20%</p> <p>Percentage of egress buffer space allocated to multicast traffic (multicast buffer partition): 75%</p>                       |

# Configuration

## IN THIS SECTION

- [CLI Quick Configuration | 724](#)
- [Configuring the Global Shared Buffer Pool for Networks with Mostly Multicast Traffic | 724](#)
- [Results | 725](#)

## CLI Quick Configuration

To quickly configure the recommended shared buffer settings for networks that carry mostly multicast traffic, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit class-of-service shared-buffer] hierarchy level:

```
[edit class-of-service shared-buffer]
set ingress percent 100
set ingress buffer-partition lossless percent 5
set ingress buffer-partition lossless-headroom percent 0
set ingress buffer-partition lossy percent 95
set egress percent 100
set egress buffer-partition lossless percent 5
set egress buffer-partition lossy percent 20
set egress buffer-partition multicast percent 75
```

## Configuring the Global Shared Buffer Pool for Networks with Mostly Multicast Traffic

### Step-by-Step Procedure

To configure the global ingress and egress shared buffer allocations and partitions for a network that carries mostly multicast traffic:

1. Configure the percentage of available (nonreserved) buffers used for the ingress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set ingress percent 100
```

2. Configure the global ingress buffer partitions for lossless, lossless-headroom, and lossy traffic:

```
[edit class-of-service shared-buffer]
user@switch# set ingress buffer-partition lossless percent 5
user@switch# set ingress buffer-partition lossless-headroom percent 0
user@switch# set ingress buffer-partition lossy percent 95
```

3. Configure the percentage of available (nonreserved) buffers used for the egress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set egress percent 100
```

4. Configure the global egress buffer partitions for lossless, lossy, and multicast queues:

```
[edit class-of-service shared-buffer]
user@switch# set egress buffer-partition lossless percent 5
user@switch# set egress buffer-partition lossy percent 20
user@switch# set egress buffer-partition multicast percent 75
```

## Results

Display the results of the configuration:

```
root@dcbg-tp-pa-02> show configuration class-of-service shared-buffer
ingress {
 percent 100;
 buffer-partition lossless {
 percent 5;
 }
 buffer-partition lossy {
 percent 95;
 }
 buffer-partition lossless-headroom {
 percent 0;
 }
}
egress {
 percent 100;
 buffer-partition lossless {
 percent 5;
 }
 buffer-partition lossy {
 percent 20;
 }
 buffer-partition multicast {
```



```
 percent 75;
 }
}
```

Verification

IN THIS SECTION

Verifying the Shared Buffer Configuration | 726

Verify that you correctly configured the shared buffer.

Verifying the Shared Buffer Configuration

Purpose

Verify that you correctly configured the ingress and egress global shared buffer pools and that you correctly partitioned the buffer among the shared buffer types.

Action

List the global shared buffer configuration using the operational mode command `show class-of-service shared-buffer`:

```
user@switch> show class-of-service shared-buffer
root@dcbg-tp-pa-02> show class-of-service shared-buffer
Ingress:
 Total Buffer : 9360.00 KB
 Dedicated Buffer : 2158.00 KB
 Shared Buffer : 7202.00 KB
 Lossless : 360.10 KB
 Lossless Headroom : 0.00 KB
 Lossy : 6841.90 KB

Lossless Headroom Utilization:
Node Device Total Used Free
0 0.00 KB 0.00 KB 0.00 KB
```

**Egress:**

```

Total Buffer : 9360.00 KB
Dedicated Buffer : 2704.00 KB
Shared Buffer : 6656.00 KB
 Lossless : 332.80 KB
 Multicast : 4992.00 KB
 Lossy : 1331.20 KB

```

**Meaning**

The `show class-of-service shared-buffer operational` command shows all of the ingress and egress global shared buffer settings, including the buffer partitioning.

For the ingress shared buffers, the command output shows:

- The total switch buffer pool is 9360 KB (9 MB).
- The dedicated buffer pool is 2158 KB. This is the size of the global ingress dedicated buffer pool when you configure the ingress shared buffer pool as 100 percent of the available (user-configurable) buffer space. This is the minimum size of the reserved, ingress dedicated ingress buffer pool (not user-configurable). If you configure the shared buffer as less than 100 percent of the available buffer pool, the remaining buffer space is added to the dedicated buffer pool.
- With the ingress shared buffer pool configured as 100 percent of the available buffers, the total size of the ingress shared buffer pool is 7202 KB.
- The ingress shared buffer pool is partitioned to allocate:
  - 360.10 KB to lossless traffic
  - No space to lossless headroom traffic
  - 6841.90 KB to lossy unicast traffic
- The Lossless Headroom Utilization field shows how much of the buffer space reserved for paused traffic is used. Because the lossless headroom buffer partition is set to 0 (zero) percent, the total amount of lossless headroom buffer space is 0 KB; therefore the amount of used and free lossless headroom buffer space is also 0 KB.

For the egress shared buffers, the command output shows:

- The total switch buffer pool is 9360 KB (9 MB).
- The dedicated buffer pool is 2704 KB. This is the size of the global egress dedicated buffer pool when you configure the egress shared buffer pool as 100 percent of the available (user-configurable) buffer space. This is the minimum size of the reserved, egress dedicated buffer pool (not user-

configurable). If you configure the shared buffer as less than 100 percent of the available buffer pool, the remaining buffer space is added to the dedicated buffer pool.

- With the egress shared buffer pool configured as 100 percent of the available buffers, the total size of the egress shared buffer pool is 6656 KB. This is less than the ingress shared buffer pool because the switch reserves more egress dedicated buffer space than ingress dedicated buffer space. (More dedicated buffer space means less shared buffer space, and more shared buffer space means less dedicated buffer space.)
- The egress shared buffer pool is partitioned to allocate:
  - 332.80 KB to lossless traffic
  - 4992 KB to multicast traffic
  - 1331.20 KB to lossy unicast traffic

**NOTE:** The output values are valid for QFX3500 and QFX3600 switches. QFX5100, EX4600, and OCX Series switches have larger buffers (12 MB instead of 9 MB), so the total buffer size and the sizes of each buffer partition are larger on those switches.

## RELATED DOCUMENTATION

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic*

*Configuring Global Ingress and Egress Shared Buffers*

*Understanding CoS Buffer Configuration*

## Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic

### IN THIS SECTION

- Requirements | 730
- Overview | 730
- Configuration | 732
- Verification | 735

Although the switch reserves some buffer space to ensure a minimum memory allocation for ports and queues, you can configure how the system uses the rest of the buffer space to optimize the buffer allocation for your particular mix of network traffic.

This example shows you the recommended configuration of the global shared buffer pool to support a network that carries mostly lossless traffic. The global shared buffer pool is memory space that all of the ports on the switch share dynamically as they need buffers. You can allocate global shared memory space to different types of buffers to better support different mixes of network traffic.



**CAUTION:** Changing the buffer configuration is a disruptive event. Traffic stops on *all* ports until buffer reprogramming is complete.

Use the default shared buffer settings (for a network with a balanced mix of lossless, best effort, and multicast traffic) or one of the recommended shared buffer configurations for your mix of network traffic (mostly best-effort unicast traffic, mostly best-effort traffic on links enabled for Ethernet PAUSE, mostly multicast traffic, or mostly lossless traffic). Either the default configuration or one of the recommended configurations provides a buffer allocation that satisfies the needs of most networks.

**NOTE:** When we discuss lossless buffers, we mean buffers that handle traffic on which you enable priority-based flow control (PFC) to ensure lossless transport. The lossless buffers are not used for best-effort traffic on a link on which you enable Ethernet PAUSE (IEEE 802.3x).

After starting from the recommended configuration, you can fine-tune the shared buffer settings, but do so with caution to prevent traffic loss due to buffer misconfiguration.

## Requirements

This example uses the following hardware and software components:

- Juniper Networks QFX3500 Switch
- Junos OS Release 12.3 or later for the QFX Series

## Overview

### IN THIS SECTION

- [Topology | 731](#)

You can configure the percentage of available (user-configurable) buffer space allocated to the global shared buffers. Any space that you do not allocate to the global shared buffer pool is added to the dedicated buffer pool. The default configuration allocates 100 percent of the available buffer space to the global shared buffers.

You can partition the ingress and egress shared buffer pools to allocate more buffers to the types of traffic your network predominantly carries, and fewer buffers to other traffic. From the buffer space allocated to the ingress shared buffer pool, you can allocate space to:

- Lossless buffers—Percentage of shared buffer pool for all lossless ingress traffic. The minimum value for the lossless buffers is 5 percent.
- Lossless headroom buffers—Percentage of shared buffer pool for packets received while a pause is asserted. If Ethernet PAUSE is configured on a port or if priority-based flow control (PFC) is configured on priorities on a port, when the port sends a pause message to the connected peer, the port uses the headroom buffers to store the packets that arrive between the time the port sends the pause message and the time the last packet arrives after the peer pauses traffic. The minimum value for the lossless headroom buffers is 0 (zero) percent. (Lossless headroom buffers are the only buffers that can have a minimum value of less than 5 percent.)

**NOTE:** On a QFX Virtual Chassis and an EX4600/EX4650 Virtual Chassis, the minimum value for the lossless headroom buffer is 3 percent.

- Lossy buffers—Percentage of shared buffer pool for all best-effort ingress traffic (best-effort unicast, multdestination, and strict-high priority traffic). The minimum value for the lossy buffers is 5 percent.

The combined percentage values of the ingress lossless, lossless headroom, and lossy buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All ingress buffer partitions must be explicitly configured, even when the lossless headroom buffer partition has a value of 0 (zero) percent.

**NOTE:** If you commit a buffer configuration for which the switch does not have sufficient resources, the switch might log an error instead of returning a commit error. In that case, a syslog message is displayed on the console. For example:

```
user@host# commit
configuration check succeeds

Message from syslogd@host at Jun 13 11:11:10 ...
host dc-pfe: Not enough Ingress Lossless headroom.(Already allocated more). Dedicated : 14340
Lossy : 47100 Lossless 4239 Headroom 21195 Avail : 20781
commit complete
```

From the buffer space allocated to the egress shared buffer pool, you can allocate space to:

- Lossless buffers—Percentage of shared buffer pool for all lossless egress queues. The minimum value for the lossless buffers is 5 percent.
- Lossy buffers—Percentage of shared buffer pool for all best-effort egress queues (best-effort unicast, and strict-high priority queues). The minimum value for the lossy buffers is 5 percent.
- Multicast buffers—Percentage of shared buffer pool for all multideestination (multicast, broadcast, and destination lookup fail) egress queues. The minimum value for the multicast buffers is 5 percent.

The combined percentage values of the egress lossless, lossy, and multicast buffer partitions must total exactly 100 percent. If the buffer percentages total more than 100 percent or less than 100 percent, the switch returns a commit error. All egress buffer partitions must be explicitly configured and must have a value of at least 5 percent.

To configure the shared buffers to support a network that carries mostly lossless traffic, more buffer space needs to be allocated to lossless buffers, and less buffer space should be allocated to lossy buffers. This example shows you how to configure the global shared buffer pool allocation that we recommend to support a network that carries mostly lossless traffic.

## Topology

[Table 135 on page 732](#) shows the configuration components for this example.

**Table 135: Components of the Recommended Shared Buffer Configuration for Lossless Network Topologies**

| Component             | Settings                                                                                                                                                                                                                                                                                                                                                                                                                                     |
|-----------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Hardware              | QFX3500 switch                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Ingress shared buffer | <p>Percentage of available ingress buffer space allocated to the ingress shared buffer: 100%</p> <p>Percentage of ingress buffer space allocated to lossless traffic (lossless buffer partition): 15%</p> <p>Percentage of ingress buffer space allocated to lossless headroom traffic (lossless headroom buffer partition): 80%</p> <p>Percentage of ingress buffer space allocated to best-effort traffic (lossy buffer partition): 5%</p> |
| Egress shared buffer  | <p>Percentage of available egress buffer space allocated to the egress shared buffer: 100%</p> <p>Percentage of egress buffer space allocated to lossless queues (lossless buffer partition): 90%</p> <p>Percentage of egress buffer space allocated to best-effort queues (lossy buffer partition): 5%</p> <p>Percentage of egress buffer space allocated to multicast traffic (multicast buffer partition): 5%</p>                         |

## Configuration

### IN THIS SECTION

- [CLI Quick Configuration | 733](#)
- [Configuring the Global Shared Buffer Pool for Networks with Mostly Lossless Traffic | 733](#)
- [Results | 734](#)

## CLI Quick Configuration

To quickly configure the recommended shared buffer settings for networks that carry mostly lossless traffic, copy the following commands, paste them in a text file, remove line breaks, change variables and details to match your network configuration, and then copy and paste the commands into the CLI at the [edit] hierarchy level:

```
[edit class-of-service shared-buffer]
set ingress percent 100
set ingress buffer-partition lossless percent 15
set ingress buffer-partition lossless-headroom percent 80
set ingress buffer-partition lossy percent 5
set egress percent 100
set egress buffer-partition lossless percent 90
set egress buffer-partition lossy percent 5
set egress buffer-partition multicast percent 5
```

## Configuring the Global Shared Buffer Pool for Networks with Mostly Lossless Traffic

### Step-by-Step Procedure

To configure the global ingress and egress shared buffer allocations and partitions for a network that carries mostly lossless traffic:

1. Configure the percentage of available (nonreserved) buffers used for the ingress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set ingress percent 100
```

2. Configure the global ingress buffer partitions for lossless, lossless-headroom, and lossy traffic:

```
[edit class-of-service shared-buffer]
user@switch# set ingress buffer-partition lossless percent 15
user@switch# set ingress buffer-partition lossless-headroom percent 80
user@switch# set ingress buffer-partition lossy percent 5
```



3. Configure the percentage of available (nonreserved) buffers used for the egress global shared buffer pool:

```
[edit class-of-service shared-buffer]
user@switch# set egress percent 100
```

4. Configure the global egress buffer partitions for lossless, lossy, and multicast queues:

```
[edit class-of-service shared-buffer]
user@switch# set egress buffer-partition lossless percent 90
user@switch# set egress buffer-partition lossy percent 5
user@switch# set egress buffer-partition multicast percent 5
```

## Results

Display the results of the configuration:

```
root@dcbg-tp-pa-02> show configuration class-of-service shared-buffer
ingress {
 percent 100;
 buffer-partition lossless {
 percent 15;
 }
 buffer-partition lossy {
 percent 5;
 }
 buffer-partition lossless-headroom {
 percent 80;
 }
}
egress {
 percent 100;
 buffer-partition lossless {
 percent 90;
 }
 buffer-partition lossy {
 percent 5;
 }
 buffer-partition multicast {
```

```
 percent 5;
 }
}
```

Verification

IN THIS SECTION

[Verifying the Shared Buffer Configuration | 735](#)

Verify that the shared buffer configuration has been created properly.

Verifying the Shared Buffer Configuration

Purpose

Verify that the ingress and egress global shared buffer pools are correctly configured and partitioned among the shared buffer types.

Action

List the global shared buffer configuration using the operational mode command `show class-of-service shared-buffer`:

```
user@switch> show class-of-service shared-buffer
root@dcbg-tp-pa-02> show class-of-service shared-buffer
Ingress:
 Total Buffer : 9360.00 KB
 Dedicated Buffer : 2158.00 KB
 Shared Buffer : 7202.00 KB
 Lossless : 1080.30 KB
 Lossless Headroom : 5761.60 KB
 Lossy : 360.10 KB

Lossless Headroom Utilization:
Node Device Total Used Free
0 5761.60 KB 0.00 KB 5761.60 KB
```

```
Egress:
 Total Buffer : 9360.00 KB
 Dedicated Buffer : 2704.00 KB
 Shared Buffer : 6656.00 KB
 Lossless : 5990.40 KB
 Multicast : 332.80 KB
 Lossy : 332.80 KB
```

## Meaning

The `show class-of-service shared-buffer operational` command shows all of the ingress and egress global shared buffer settings, including the buffer partitioning.

For the ingress shared buffers, the command output shows:

- The total switch buffer pool is 9360 KB (9 MB).
- The dedicated buffer pool is 2158 KB. This is the size of the global ingress dedicated buffer pool when you configure the ingress shared buffer pool as 100 percent of the available (user-configurable) buffer space. This is the minimum size of the reserved, ingress dedicated ingress buffer pool (not user-configurable). If you configure the shared buffer as less than 100 percent of the available buffer pool, the remaining buffer space is added to the dedicated buffer pool.
- With the ingress shared buffer pool configured as 100 percent of the available buffers, the total size of the ingress shared buffer pool is 7202 KB.
- The ingress shared buffer pool is partitioned to allocate:
  - 1080 KB to lossless traffic
  - 5761.60 KB to lossless headroom traffic
  - 360.10 KB to lossy unicast traffic
- The Lossless Headroom Utilization field shows how much of the buffer space reserved for paused traffic is used. Of the total available lossless headroom buffer space of 5761.60 KB, currently no buffer space is being used, so all 5761.60 KB of buffer space is free.

For the egress shared buffers, the command output shows:

- The total switch buffer pool is 9360 KB (9 MB).
- The dedicated buffer pool is 2704 KB. This is the size of the global egress dedicated buffer pool when you configure the egress shared buffer pool as 100 percent of the available (user-configurable) buffer space. This is the minimum size of the reserved, egress dedicated buffer pool (not user-

configurable). If you configure the shared buffer as less than 100 percent of the available buffer pool, the remaining buffer space is added to the dedicated buffer pool.

- With the egress shared buffer pool configured as 100 percent of the available buffers, the total size of the egress shared buffer pool is 6656 KB. This is less than the ingress shared buffer pool because the switch reserves more egress dedicated buffer space than ingress dedicated buffer space. (More dedicated buffer space means less shared buffer space, and more shared buffer space means less dedicated buffer space.)
- The egress shared buffer pool is partitioned to allocate:
  - 5990.40 KB to lossless traffic
  - 332.80 KB to multicast traffic
  - 332.80 KB to lossy unicast traffic

**NOTE:** The output values are valid for QFX3500 and QFX3600 switches. QFX5100 and EX4600 switches have larger buffers (12MB instead of 9MB), so the total buffer size and the sizes of each buffer partition are larger on QFX5100 and EX4600 switches.

## RELATED DOCUMENTATION

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled*

*Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic*

*Configuring Global Ingress and Egress Shared Buffers*

*Understanding CoS Buffer Configuration*



## Learn About Technology

---

Learn About Technology | 739

---

## CHAPTER 9

# Learn About Technology

**IN THIS CHAPTER**

- [Data Center Technology Overview Videos | 739](#)

## Data Center Technology Overview Videos

**IN THIS SECTION**

- [Learn About Video: Why Do We Need an IP Fabric? | 739](#)
- [Learn About Video: What is the Best Control Plane Protocol to Use in a Data Center IP Fabric? | 740](#)
- [Learn About Video: Why Use an Overlay Network in a Data Center? | 740](#)
- [Conceptual Documents That Contain Technology Overview Videos | 740](#)

Juniper Information Experience (iX) videos provide brief, high-level overviews of data center technologies and concepts. Each video runs approximately one-and-a-half to two minutes in length. This document contains SDN-related videos and links to conceptual documents that contain other data center technology videos:

### Learn About Video: Why Do We Need an IP Fabric?

The video *Why Do We Need an IP Fabric?* presents a brief overview of IP Fabric use cases.



Video: [Why Do We Need an IP Fabric?](#)

## Learn About Video: What is the Best Control Plane Protocol to Use in a Data Center IP Fabric?

The video *What is the Best Control Plane Protocol to Use in a Data Center IP Fabric?* presents a brief overview of the arguments for using Border Gateway Protocol (BGP) as the data center IP fabric control plane protocol.



Video: [What is the Best Control Plane Protocol to Use in a Data Center IP Fabric?](#)

## Learn About Video: Why Use an Overlay Network in a Data Center?

The video *Why Use an Overlay Network in a Data Center?* presents a brief overview of the advantages of data center overlay networks.



Video: [Why Use an Overlay Network in a Data Center?](#)

## Conceptual Documents That Contain Technology Overview Videos

The following conceptual documents include brief video overviews of the technology:

- [Understanding DCB Features and Requirements](#)
- [Understanding CoS Hierarchical Port Scheduling \(ETS\)](#)
- [Understanding CoS Flow Control \(Ethernet PAUSE and PFC\)](#)
- [Understanding DCBX](#)
- [Understanding PFC Functionality Across Layer 3 Interfaces](#)
- [Virtual Chassis Fabric Overview](#)
- [Understanding In-Service Software Upgrade \(ISSU\) and In-Service Software Upgrade \(ISSU\) System Requirements](#) (same video)

# 7

PART

## Configuration Statements and Operational Commands

---

[Monitoring Interfaces That Have CoS Components | 742](#)

[Monitoring CoS Classifiers | 744](#)

[Monitoring CoS Forwarding Classes | 746](#)

[Monitoring CoS Rewrite Rules | 750](#)

[Monitoring CoS Code-Point Value Aliases | 752](#)

[Monitoring CoS Scheduler Maps | 754](#)

[Junos CLI Reference Overview | 756](#)

---



# Monitoring Interfaces That Have CoS Components

IN THIS SECTION

- Purpose | 742
- Action | 742
- Meaning | 742

## Purpose

Use the monitoring functionality to display details about the physical and logical interfaces and the CoS components assigned to them.

## Action

To monitor interfaces that have CoS components in the CLI, enter the command:

```
user@switch> show class-of-service interface
```

To monitor a specific interface in the CLI, enter the command:

```
user@switch> show class-of-service interface interface-name
```

## Meaning

[Table 136 on page 742](#) summarizes key output fields for CoS interfaces.

**Table 136: Summary of Key CoS Interfaces Output Fields**

| Field              | Values                                                             |
|--------------------|--------------------------------------------------------------------|
| Physical interface | Name of a physical interface to which CoS components are assigned. |

Table 136: Summary of Key CoS Interfaces Output Fields *(Continued)*

| Field                         | Values                                                                                                                                                                                                                                        |
|-------------------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Index                         | Index of this interface or the internal index of a specific object.                                                                                                                                                                           |
| Queues supported              | Number of queues you can configure on the interface.                                                                                                                                                                                          |
| Queues in use                 | Number of queues currently configured.                                                                                                                                                                                                        |
| Scheduler map                 | Name of the scheduler map associated with this interface.                                                                                                                                                                                     |
| Congestion-notification       | Status of congestion notification (enabled or disabled).<br><br><b>NOTE:</b> OCX Series switches do not support congestion notification profiles.                                                                                             |
| Rewrite Input IEEE Code-point | (Fibre Channel NP_Port interfaces only) IEEE 802.1p code point (priority) the interface assigns to incoming Fibre Channel (FC) traffic when the interface encapsulates the FC traffic in Ethernet before forwarding it onto the FCoE network. |
| Logical Interface             | Name of a logical interface on the physical interface to which CoS components are assigned.                                                                                                                                                   |
| Object                        | Category of an object—for example, classifier, scheduler-map, or rewrite.                                                                                                                                                                     |
| Name                          | Name of the object—for example, ba-classifier.                                                                                                                                                                                                |
| Type                          | Type of the object—for example, ieee8021p for a classifier.                                                                                                                                                                                   |

## RELATED DOCUMENTATION

*Assigning CoS Components to Interfaces*

# Monitoring CoS Classifiers

IN THIS SECTION

- Purpose | 744
- Action | 744
- Meaning | 744

## Purpose

Display the mapping of incoming CoS values to forwarding class and loss priority for each classifier.

## Action

To monitor CoS classifiers in the CLI, enter the CLI command:

```
user@switch> show class-of-service classifier
```

To monitor a particular classifier in the CLI, enter the CLI command:

```
user@switch> show class-of-service classifier name classifier-name
```

To monitor a particular type of classifier in the CLI, enter the CLI command:

```
user@switch> show class-of-service classifier type classifier-type
```

## Meaning

[Table 137 on page 745](#) summarizes key output fields for CoS classifiers.

Table 137: Summary of Key CoS Classifier Output Fields

| Field            | Values                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                |
|------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Classifier       | Name of a classifier.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
| Code point type  | <p>Type of classifier:</p> <ul style="list-style-type: none"> <li>dscp—All classifiers of the DSCP type.</li> <li>ieee-802.1—All classifiers of the IEEE 802.1 type.</li> <li>ieee-mcast—All classifiers of the IEEE 802.1 multicast type.</li> </ul> <p><b>NOTE:</b> QFX10000 switches do not use different classifiers for unicast and multideestination (multicast, broadcast, destination lookup fail) traffic, so multicast-specific classifiers are not supported.</p> <ul style="list-style-type: none"> <li>exp—All classifiers of the MPLS exp type.</li> </ul> <p><b>NOTE:</b> OCX Series switches do not support MPLS.</p> |
| Index            | Internal index of the classifier.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                     |
| Code point       | DSCP or IEEE 802.1 code point value of the incoming packets, in bits. These values are used for classification.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |
| Forwarding Class | Name of the forwarding class that the classifier assigns to an incoming packet. This class affects the forwarding and scheduling policies that are applied to the packet as it transits the switch.                                                                                                                                                                                                                                                                                                                                                                                                                                   |
| Loss Priority    | Loss priority value that the classifier assigns to the incoming packet based on its code point value.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |

# Monitoring CoS Forwarding Classes

## IN THIS SECTION

- Purpose | 746
- Action | 746
- Meaning | 746

## Purpose

Use the monitoring functionality to view the current assignment of CoS forwarding classes to queue numbers on the system.

## Action

To monitor CoS forwarding classes in the CLI, enter the following CLI command:

```
user@switch> show class-of-service forwarding-class
```

## Meaning

Some switches use different forwarding classes, output queues, and classifiers for unicast and multideestination (multicast, broadcast, destination lookup fail) traffic. These switches support 12 forwarding classes and output queues, eight for unicast traffic and four for multideestination traffic.

Some switches use the same forwarding classes, output queues, and classifiers for unicast and multideestination traffic. These switches support eight forwarding classes and eight output queues.

[Table 138 on page 747](#) summarizes key output fields on switches that use different forwarding classes and output queues for unicast and multideestination traffic.

**Table 138: Summary of Key CoS Forwarding Class Output Fields on Switches that Separate Unicast and Multidestination Traffic**

| Field            | Values                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |
|------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Forwarding Class | <p>Names of forwarding classes assigned to queue numbers. By default, the following unicast forwarding classes are assigned to queues 0, 3, 4, and 7, respectively:</p> <ul style="list-style-type: none"> <li>• best-effort—Provides no special CoS handling of packets. Loss priority is typically not carried in a CoS value.</li> <li>• fcoe—Provides guaranteed delivery for Fibre Channel over Ethernet (FCoE) traffic.</li> <li>• no-loss—Provides guaranteed delivery for TCP lossless traffic</li> <li>• network-control—Packets can be delayed but not dropped.</li> </ul> <p>By default, the following multidestination forwarding class is assigned to queue 8:</p> <ul style="list-style-type: none"> <li>• mcast—Provides no special CoS handling of packets.</li> </ul> |
| Queue            | <p>Queue number corresponding to (mapped to) the forwarding class name.</p> <p>By default, four queues (0, 3, 4, and 7) are assigned to unicast forwarding classes and one queue (8) is assigned to a multidestination forwarding class:</p> <ul style="list-style-type: none"> <li>• Queue 0—best-effort</li> <li>• Queue 3—fcoe</li> <li>• Queue 4—no-loss</li> <li>• Queue 7—network-control</li> <li>• Queue 8—mcast</li> </ul>                                                                                                                                                                                                                                                                                                                                                    |

**Table 138: Summary of Key CoS Forwarding Class Output Fields on Switches that Separate Unicast and Multidestination Traffic *(Continued)***

| Field   | Values                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
|---------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| No-Loss | <p>Packet drop attribute associated with each forwarding class:</p> <ul style="list-style-type: none"> <li>• Disabled—The forwarding class is configured for lossy transport (packets might drop during periods of congestion)</li> <li>• Enabled—The forwarding class is configured for lossless transport</li> </ul> <p><b>NOTE:</b> To achieve lossless transport, you must ensure that priority-based flow control (PFC) and DCBX are properly configured on the lossless priority (IEEE 802.1p code point), and that sufficient port bandwidth is reserved for the lossless traffic flows.</p> <p>OCX Series switches do not support lossless transport.</p> |

**NOTE:** OCX Series switches do not support the default lossless forwarding classes `fcoe` and `no-loss`, and do not support the no-loss packet drop attribute used to configure lossless forwarding classes. On OCX Series switches, do not map traffic to the default `fcoe` and `no-loss` forwarding classes (both of these default forwarding classes carry the no-loss packet drop attribute), and do not configure the no-loss packet drop attribute on forwarding classes.

Table 139 on page 749 summarizes key output fields on switches that use the same forwarding classes and output queues for unicast and multidestination traffic.

**Table 139: Summary of Key CoS Forwarding Class Output Fields on Switches That Do Not Separate Unicast and Multidestination Traffic**

| Field            | Values                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |
|------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Forwarding Class | <p>Names of forwarding classes assigned to queue numbers. By default, the following forwarding classes are assigned to queues 0, 3, 4, and 7, respectively:</p> <ul style="list-style-type: none"> <li>• best-effort—Provides no special CoS handling of packets. Loss priority is typically not carried in a CoS value.</li> <li>• fcoe—Provides guaranteed delivery for Fibre Channel over Ethernet (FCoE) traffic.</li> <li>• no-loss—Provides guaranteed delivery for TCP lossless traffic</li> <li>• network-control—Packets can be delayed but not dropped.</li> </ul> |
| Queue            | <p>Queue number corresponding to (mapped to) the forwarding class name.</p> <p>By default, four queues (0, 3, 4, and 7) are assigned to forwarding classes:</p> <ul style="list-style-type: none"> <li>• Queue 0—best-effort</li> <li>• Queue 3—fcoe</li> <li>• Queue 4—no-loss</li> <li>• Queue 7—network-control</li> </ul>                                                                                                                                                                                                                                                |



Table 139: Summary of Key CoS Forwarding Class Output Fields on Switches That Do Not Separate Unicast and Multidestination Traffic *(Continued)*

| Field   | Values                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                           |
|---------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| No-Loss | <p>Packet drop attribute associated with each forwarding class:</p> <ul style="list-style-type: none"><li>• Disabled—The forwarding class is configured for lossy transport (packets might drop during periods of congestion).</li><li>• Enabled—The forwarding class is configured for lossless transport.</li></ul> <p><b>NOTE:</b> To achieve lossless transport, you must ensure that priority-based flow control (PFC) and DCBX are properly configured on the lossless priority (IEEE 802.1p code point), and that sufficient port bandwidth is reserved for the lossless traffic flows.</p> <p>OCX Series switches do not support lossless transport.</p> |

# Monitoring CoS Rewrite Rules

IN THIS SECTION

- Purpose | 750
- Action | 751
- Meaning | 751

## Purpose

Use the monitoring functionality to display information about CoS value rewrite rules, which are based on the forwarding class and loss priority.

## Action

To monitor CoS rewrite rules in the CLI, enter the CLI command:

```
user@switch> show class-of-service rewrite-rule
```

To monitor a particular rewrite rule in the CLI, enter the CLI command:

```
user@switch> show class-of-service rewrite-rule name rewrite-rule-name
```

To monitor a particular type of rewrite rule (for example, DSCP, DSCP IPv6, IEEE-802.1, or MPLS EXP) in the CLI, enter the CLI command:

```
user@switch> show class-of-service rewrite-rule type rewrite-rule-type
```

## Meaning

[Table 140 on page 751](#) summarizes key output fields for CoS rewrite rules.

**Table 140: Summary of Key CoS Rewrite Rule Output Fields**

| Field           | Values                                                                                                                                                                                                                                                                                      |
|-----------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Rewrite rule    | Name of the rewrite rule.                                                                                                                                                                                                                                                                   |
| Code point type | <div>Rewrite rule type:<ul style="list-style-type: none"><li>dscp—For IPv4 DiffServ traffic.</li><li>dscp-ipv6—For IPv6 Diffserv traffic.</li><li>ieee-802.1—For Layer 2 traffic.</li><li>exp—For MPLS traffic.</li></ul><p><b>NOTE:</b> OCX Series switches do not support MPLS.</p></div> |
| Index           | Internal index for the rewrite rule.                                                                                                                                                                                                                                                        |

Table 140: Summary of Key CoS Rewrite Rule Output Fields *(Continued)*

| Field            | Values                                                                                                                                                                                                                                            |
|------------------|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Forwarding class | <p>Name of the forwarding class that is used to determine CoS values for rewriting in combination with loss priority.</p> <p>Rewrite rules are applied to CoS values in outgoing packets based on forwarding class and loss priority setting.</p> |
| Loss priority    | <p>Level of loss priority that is used to determine CoS values for rewriting in combination with forwarding class.</p>                                                                                                                            |
| Code point       | <p>Rewrite code point value.</p>                                                                                                                                                                                                                  |

RELATED DOCUMENTATION

| [Defining CoS Rewrite Rules](#)

# Monitoring CoS Code-Point Value Aliases

IN THIS SECTION

- [Purpose | 752](#)
- [Action | 753](#)
- [Meaning | 753](#)

## Purpose

Use the monitoring functionality to display information about the CoS code-point value aliases that the system is currently using to represent DSCP and IEEE 802.1p code point bits.

## Action

To monitor CoS value aliases in the CLI, enter the CLI command:

```
user@switch> show class-of-service code-point-aliases
```

To monitor a specific type of code-point alias (DSCP, DSCP IPv6, IEEE 802.1, or MPLS EXP) in the CLI, enter the CLI command:

```
user@switch> show class-of-service code-point-aliases ieee-802.1
```

## Meaning

[Table 141 on page 753](#) summarizes key output fields for CoS value aliases.

**Table 141: Summary of Key CoS Value Alias Output Fields**

| Field           | Values                                                                                                                                                                                                                                                                                                                                                                                                                                             |
|-----------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Code point type | Type of the CoS value: <ul style="list-style-type: none"><li>• dscp—Examines Layer 3 packet headers for IP packet classification.</li><li>• dscp-ipv6—Examines Layer 3 packet headers for IPv6 packet classification.</li><li>• ieee-802.1—Examines Layer 2 packet headers for packet classification.</li><li>• exp—Examines MPLS packet headers for packet classification.</li></ul> <p><b>NOTE:</b> OCX Series switches do not support MPLS.</p> |
| Alias           | Name given to a set of bits—for example, af11 is a name for bits 001010.                                                                                                                                                                                                                                                                                                                                                                           |
| Bit pattern     | Set of bits associated with the alias.                                                                                                                                                                                                                                                                                                                                                                                                             |

### RELATED DOCUMENTATION

*Defining CoS Code-Point Aliases*

# Monitoring CoS Scheduler Maps

IN THIS SECTION

- Purpose | 754
- Action | 754
- Meaning | 754

## Purpose

Use the monitoring functionality to display assignments of CoS forwarding classes to schedulers.

## Action

To monitor CoS scheduler maps in the CLI, enter the CLI command:

```
user@switch> show class-of-service scheduler-map
```

To monitor a specific scheduler map in the CLI, enter the CLI command:

```
user@switch> show class-of-service scheduler-map scheduler-map-name
```

## Meaning

[Table 142 on page 754](#) summarizes key output fields for CoS scheduler maps.

**Table 142: Summary of Key CoS Scheduler Maps Output Fields**

| Field         | Values                                                              |
|---------------|---------------------------------------------------------------------|
| Scheduler map | Name of a scheduler map that maps forwarding classes to schedulers. |

Table 142: Summary of Key CoS Scheduler Maps Output Fields *(Continued)*

| Field            | Values                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                              |
|------------------|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Index            | Index of a specific object—scheduler maps, schedulers, or drop profiles.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
| Scheduler        | Name of a scheduler that controls queue properties such as bandwidth and scheduling priority.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                       |
| Forwarding class | Name(s) of the forwarding class(es) to which the scheduler is mapped.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                               |
| Transmit rate    | Guaranteed minimum bandwidth configured on the queue mapped to the scheduler. On strict-high priority queues on QFX10000 switches, defines the maximum amount of traffic on the queue that is treated as strict-high priority traffic.                                                                                                                                                                                                                                                                                                                                                              |
| Priority         | <p>Scheduling priority of traffic on a queue:</p> <ul style="list-style-type: none"> <li>strict-high or high—Packets on a strict-high priority queue are transmitted first, before all other traffic, up to the configured maximum bandwidth (shaping rate). On QFX3500, QFX3600, EX4600, and OCX series switches, and on QFabric system, only one queue can be configured as strict-high or high priority. On QFX10000 switches, you can configure more than one strict-high priority queue.</li> <li>low—Packets in this queue are transmitted after packets in the strict-high queue.</li> </ul> |
| Drop Profiles    | Name and index of a drop profile that is mapped to a specific loss priority and protocol pair. The drop profile determines the way best effort queues drop packets during periods of congestion.                                                                                                                                                                                                                                                                                                                                                                                                    |
| Loss Priority    | Packet loss priority mapped to the drop profile. You can configure different drop profiles for low, medium-high, and high loss priority traffic.                                                                                                                                                                                                                                                                                                                                                                                                                                                    |
| Protocol         | Transport protocol of the drop profile for the particular priority.                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                                 |

Table 142: Summary of Key CoS Scheduler Maps Output Fields *(Continued)*

| Field | Values                    |
|-------|---------------------------|
| Name  | Name of the drop profile. |

# Junos CLI Reference Overview

We've consolidated all Junos CLI commands and configuration statements in one place. Learn about the syntax and options that make up the statements and commands and understand the contexts in which you'll use these CLI elements in your network configurations and operations.

- [Junos CLI Reference](#)

Click the links to access Junos OS and Junos OS Evolved configuration statement and command summary topics.

- [Configuration Statements](#)
- [CLI Commands](#)