

Traffic Management User Guide (QFX Series Switches and EX4600 Switches)

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 (QFX Series Switches and EX4600 Switches)
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

1

About This Guide | xiii

Basic CoS Configuration

CoS Overview | 2

Overview of Junos OS CoS | 2

Overview of Policers | 5

Configuring CoS | 14

Understanding Junos CoS Components | 21

Understanding CoS Packet Flow | 26

Understanding Default CoS Settings | 30

CoS Support on QFX Series Switches, EX4600 Line of Switches, and QFabric Systems | 44

CoS on Interfaces | 60

CoS Inputs and Outputs Overview | 60

CoS on Virtual Chassis Switch Ports | 61

CoS on Virtual Chassis Fabric (VCF) EX4300 Leaf Devices (Mixed Mode) | 66

Understanding CoS on OVSDB-Managed VXLAN Interfaces | 73

Configuring CoS on OVSDB-Managed VXLAN Interfaces | 78

Assigning CoS Components to Interfaces | 87

CoS Code-Point Aliases | 90

Understanding CoS Code-Point Aliases | 90

Defining CoS Code-Point Aliases | 93

Monitoring CoS Code-Point Value Aliases | 94

CoS Classifiers | 96

Understanding CoS Classifiers | 96

Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p) | 106

Example: Configuring Classifiers | 108

Requirements | 110

Overview | 110

Verification | 111

Example: Configuring Unicast Classifiers | 113

Requirements | 114

Overview | 114

Verification | 115

Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers | 117

Requirements | 118

Overview | 119

Verification | 120

Understanding Host Inbound Traffic Classification | 121

Configuring a Global MPLS EXP Classifier | 122

Monitoring CoS Classifiers | 123

CoS Rewrite Rules | 125

Understanding CoS Rewrite Rules | 125

Defining CoS Rewrite Rules | 128

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130

Troubleshooting an Unexpected Rewrite Value | 145

Understanding CoS MPLS EXP Classifiers and Rewrite Rules | 147

Configuring Rewrite Rules for MPLS EXP Classifiers | 151

Monitoring CoS Rewrite Rules | 153

CoS Forwarding Classes and Forwarding Class Sets | 155

Understanding CoS Forwarding Classes | 155

Defining CoS Forwarding Classes | 162

Forwarding Policy Options Overview | 164

Configuring CoS-Based Forwarding | 166

Example: Configuring CoS-Based Forwarding | 170

Example: Configuring Forwarding Classes | 174

Requirements | 175

Overview | 175

Example 1: Configuring Forwarding Classes for Switches Except QFX10000 | 177

Verification | 178

Example 2: Configuring Forwarding Classes for QFX10000 Switches | 179

Verification | 180

Understanding CoS Forwarding Class Sets (Priority Groups) | 181

Defining CoS Forwarding Class Sets | 183

Example: Configuring Forwarding Class Sets | 184

Requirements | 185

Overview | 185

Verification | 187

Monitoring CoS Forwarding Classes | 189

Lossless Traffic Flows, Ethernet PAUSE Flow Control, and PFC | 194

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 194

Configuring CoS PFC (Congestion Notification Profiles) | 216

Understanding CoS Flow Control (Ethernet PAUSE and PFC) | 220

Enabling and Disabling CoS Symmetric Ethernet PAUSE Flow Control | 233

Configuring CoS Asymmetric Ethernet PAUSE Flow Control | 234

Understanding PFC Functionality Across Layer 3 Interfaces | 236

Example: Configuring PFC Across Layer 3 Interfaces | 239

Requirements | 240

Overview | 240

Configuration | 246

Verification | 258

Understanding PFC Using DSCP at Layer 3 for Untagged Traffic | 266

Configuring DSCP-based PFC for Layer 3 Untagged Traffic | 268

CoS and Host Outbound Traffic | 271

Understanding Host Routing Engine Outbound Traffic Queues and Defaults | 271

Changing the Host Outbound Traffic Default Queue Mapping | 274

Weighted Random Early Detection (WRED) and Explicit Congestion Notification (ECN)

WRED and Drop Profiles | 276

Understanding CoS WRED Drop Profiles | 276

Configuring CoS WRED Drop Profiles | 284

Drop Profiles on Switches Except QFX10000 | 285

Drop Profiles on QFX 10000 Switches | 285

Example: Configuring WRED Drop Profiles | 286

Requirements | 287

Overview | 287

Configuring WRED Drop Profiles on Switches Except QFX10000 | 288

Verification | 290

Configuring WRED Drop Profiles on QFX10000 Switches | 291

Verification | 292

Configuring CoS Drop Profile Maps | 293

Example: Configuring Drop Profile Maps | 293

Requirements | 295

Overview | 295

Verification | 295

Explicit Congestion Notification (ECN) | 297

Understanding CoS Explicit Congestion Notification | 297

Example: Configuring ECN | 307

Requirements | 307

Overview | 307

Configuration | 310

Verification | 313

Data Center Quantized Congestion Notification (DCQCN) | 314

Understanding Data Center Quantized Congestion Notification (DCQCN) | 315

| Configuring Data Center Quantized Congestion Notification (DCQCN) | 316

CoS Queue Schedulers, Traffic Control Profiles, and Hierarchical Port Scheduling (ETS)

Queue Schedulers and Scheduling Priority | 321

Understanding Default CoS Scheduling and Classification | 321

Understanding CoS Scheduling Behavior and Configuration Considerations | 332

Understanding CoS Output Queue Schedulers | 338

Defining CoS Queue Schedulers | 346

Example: Configuring Queue Schedulers | 350

| Requirements | 352

| Overview | 352

| Verification | 356

Defining CoS Queue Scheduling Priority | 358

Example: Configuring Queue Scheduling Priority | 360

| Requirements | 361

| Overview | 361

| Verification | 363

Monitoring CoS Scheduler Maps | 365

Port Scheduling and Shaping | 368

Understanding CoS Port Schedulers | 368

Defining CoS Queue Schedulers for Port Scheduling | 382

Example: Configuring Queue Schedulers for Port Scheduling | 386

| Requirements | 388

| Overview | 388

| Verification | 391

CoS Port Shaping | 393

| Understanding Port Shaping | 393

| Configuring Port Shaping | 393

Troubleshooting Egress Bandwidth Issues | 396

Troubleshooting Egress Bandwidth That Exceeds the Configured Minimum Bandwidth | 396

Troubleshooting Egress Bandwidth That Exceeds the Configured Maximum Bandwidth | 398

Troubleshooting Egress Queue Bandwidth Impacted by Congestion | 399

Traffic Control Profiles and Priority Group Scheduling | 401

Understanding CoS Traffic Control Profiles | 401

Understanding CoS Priority Group Scheduling | 403

Understanding CoS Virtual Output Queues (VOQs) | 406

Defining CoS Traffic Control Profiles (Priority Group Scheduling) | 412

Example: Configuring Traffic Control Profiles (Priority Group Scheduling) | 414

Requirements | 415

Overview | 415

Verification | 416

Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth | 417

Example: Configuring Minimum Guaranteed Output Bandwidth | 421

Requirements | 423

Overview | 423

Verification | 425

Understanding CoS Priority Group Shaping and Queue Shaping (Maximum Bandwidth) | 428

Example: Configuring Maximum Output Bandwidth | 431

Requirements | 433

Overview | 433

Verification | 434

Hierarchical Port Scheduling (ETS) | 438

Understanding CoS Hierarchical Port Scheduling (ETS) | 438

Example: Configuring CoS Hierarchical Port Scheduling (ETS) | 445

Requirements | 446

Overview | 446

Configuration | 452

Verification | 466

Disabling the ETS Recommendation TLV | 480

Data Center Bridging and Lossless FCoE

Data Center Bridging | 482

Understanding DCB Features and Requirements | 482

Understanding DCBX | 486

Configuring the DCBX Mode | 496

Configuring DCBX Autonegotiation | 497

Understanding DCBX Application Protocol TLV Exchange | 500

Defining an Application for DCBX Application Protocol TLV Exchange | 504

Configuring an Application Map for DCBX Application Protocol TLV Exchange | 506

Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange | 507

Example: Configuring DCBX Application Protocol TLV Exchange | 509

Requirements | 510

Overview | 510

Configuration | 515

Verification | 518

Lossless FCoE | 524

Example: Configuring CoS PFC for FCoE Traffic | 524

Requirements | 525

Overview | 525

Configuration | 528

Verification | 535

Example: Configuring CoS for FCoE Transit Switch Traffic Across an MC-LAG | 538

Requirements | 539

Overview | 539

Configuration | 546

Verification | 559

Example: Configuring CoS Using ELS for FCoE Transit Switch Traffic Across an MC-LAG | 572

Requirements | 573

Overview | 573

Configuration | **580**

Verification | **595**

Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch) | **608**

Requirements | **608**

Overview | **608**

Configuration | **611**

Verification | **614**

Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface | **620**

Requirements | **620**

Overview | **621**

Configuration | **624**

Verification | **627**

Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces | **633**

Requirements | **633**

Overview | **634**

Configuration | **639**

Verification | **644**

Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI) | **652**

Requirements | **653**

Overview | **653**

Configuration | **660**

Verification | **668**

Troubleshooting Dropped FCoE Traffic | **678**

5

CoS Buffers and the Shared Buffer Pool

CoS Buffers Overview | **684**

Understanding CoS Buffer Configuration | **684**

Configuring Global Ingress and Egress Shared Buffers | **708**

Configuring Ingress and Egress Dedicated Buffers | **710**

- Decreasing the Global Dedicated Buffer | 711
- Configuring and Applying Dedicated Buffer Profiles | 713

Shared Buffer Pool Examples | 717

Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic | 717

- Requirements | 718
- Overview | 718
- Configuration | 720
- Verification | 723

Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled | 726

- Requirements | 727
- Overview | 727
- Configuration | 729
- Verification | 732

Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 734

- Requirements | 735
- Overview | 735
- Configuration | 737
- Verification | 740

Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 743

- Requirements | 744
- Overview | 744
- Configuration | 746
- Verification | 749

CoS on EVPN VXLANs

CoS Support on EVPN VXLANs | 753

Understanding CoS on VXLAN Interfaces | 753

Configuring CoS on VXLAN Interfaces | 754

Implementing CoS on VXLAN Interfaces (Junos OS Evolved) | 757

CoS Limitations on VXLANs | 759

Configuration Statements and Operational Commands

Junos CLI Reference Overview | 761

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. This guide applies to all QFX Series and the EX4600 line of switches.

1

PART

Basic CoS Configuration

[CoS Overview](#) | 2

[CoS on Interfaces](#) | 60

[CoS Code-Point Aliases](#) | 90

[CoS Classifiers](#) | 96

[CoS Rewrite Rules](#) | 125

[CoS Forwarding Classes and Forwarding Class Sets](#) | 155

[Lossless Traffic Flows, Ethernet PAUSE Flow Control, and PFC](#) | 194

[CoS and Host Outbound Traffic](#) | 271

CHAPTER 1

CoS Overview

IN THIS CHAPTER

- Overview of Junos OS CoS | 2
- Overview of Policers | 5
- Configuring CoS | 14
- Understanding Junos CoS Components | 21
- Understanding CoS Packet Flow | 26
- Understanding Default CoS Settings | 30
- CoS Support on QFX Series Switches, EX4600 Line of Switches, and QFabric Systems | 44

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

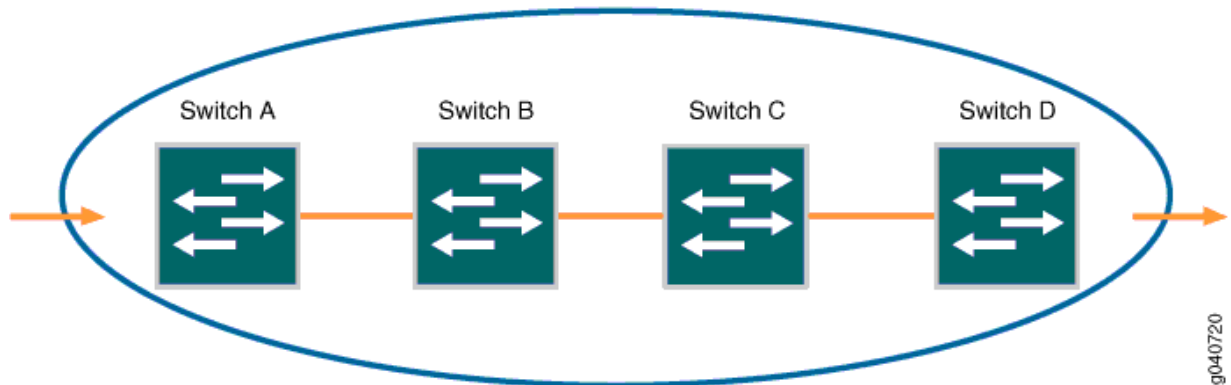
Junos OS CoS works by examining traffic entering the edge of your network. The switch classifies traffic into defined service groups to provide the special treatment of traffic across the network. For example, you can send voice traffic across certain links and data traffic across other links. In addition, 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 | 21](#)

[Understanding CoS Packet Flow | 26](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

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).

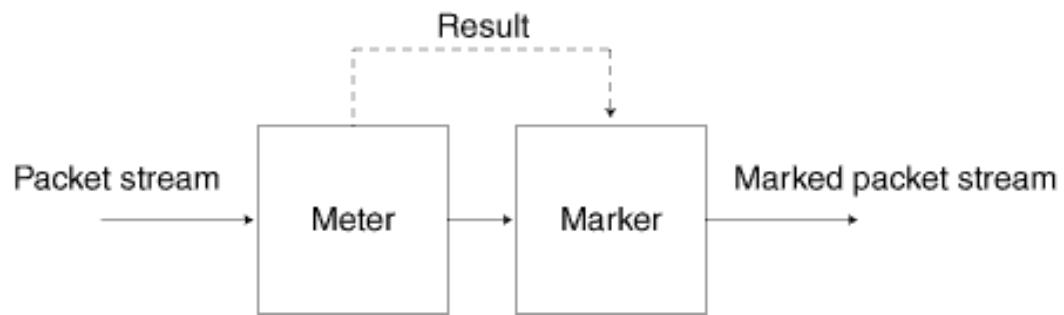
As an example:

- 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

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 16](#) 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" on page 93 (QFX10000 only) "Example: Configuring Classifiers" on page 108 (Except QFX10000) "Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p)" on page 106 (Except NFX250 and QFX10000) "Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers" on page 117 "Changing the Host Outbound Traffic Default Queue Mapping" on page 274 "Example: Configuring Forwarding Classes" on page 174 "Defining CoS Rewrite Rules" on page 128 (Except NFX250) "Enabling and Disabling CoS Symmetric Ethernet PAUSE Flow Control" on page 233 (Except NFX250 and OCX1100) "Configuring CoS Asymmetric Ethernet PAUSE Flow Control" on page 234 "Assigning CoS Components to Interfaces" on page 87

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" on page 286 "Example: Configuring Drop Profile Maps" on page 293 <i>Example: Configuring ECN</i>

Table 2: CoS Configuration Tasks *(Continued)*

CoS Configuration Task	Platforms Supported	Links
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.	<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" on page 350 • "Example: Configuring Queue Scheduling Priority" on page 360 • (QFX10000 only) "Example: Configuring Queue Schedulers for Port Scheduling" on page 386
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.	<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)" on page 412 • (Except NFX250) "Example: Configuring Traffic Control Profiles (Priority Group Scheduling)" on page 414 • "Example: Configuring Minimum Guaranteed Output Bandwidth" on page 421 • (Except NFX250) "Example: Configuring Maximum Output Bandwidth" on page 431

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" on page 184 • "Example: Configuring CoS Hierarchical Port Scheduling (ETS)" on page 445 • (Except OCX1100) "Disabling the ETS Recommendation TLV" on page 480
<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" on page 509 • "Configuring the DCBX Mode" on page 496 • "Configuring DCBX Autonegotiation" on page 497 • "Defining an Application for DCBX Application Protocol TLV Exchange" on page 504 • "Configuring an Application Map for DCBX Application Protocol TLV Exchange" on page 506 • "Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange" on page 507

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" on page 524 Example: Configuring CoS for FCoE Transit Switch Traffic Across an MC-LAG "Configuring CoS PFC (Congestion Notification Profiles)" on page 216 (QFX3500 and QFabric only) <i>Example: Configuring IEEE 802.1p Priority Remapping on an FCoE-FC Gateway</i> "Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces" on page 633 "Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic (FCoE Transit Switch)" on page 608 "Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface" on page 620 (QFX3500, NFX250, and QFabric only) <i>Configuring CoS Fixed Classifier Rewrite Values for Native FC Interfaces (NP_Ports)</i> "Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications (FCoE and iSCSI)" on page 652

Understanding Junos CoS Components

IN THIS SECTION

- [Code-Point Aliases | 21](#)
- [Policers | 21](#)
- [Classifiers | 21](#)
- [Forwarding Classes | 22](#)
- [Forwarding Class Sets | 23](#)
- [Flow Control \(Ethernet PAUSE, PFC, and ECN\) | 24](#)
- [WRED Profiles and Tail Drop | 25](#)
- [Schedulers | 25](#)
- [Rewrite Rules | 26](#)

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 multdestination 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 multdestination forwarding classes map to multdestination output queues. You cannot assign unicast traffic and multdestination traffic to the same forwarding class or to the same output queue. Switches that require the separation of unicast and multdestination traffic have 12 output queues, 8 for unicast traffic and 4 for multdestination traffic.

On switches that do not separate unicast and multdestination traffic, unicast and multdestination 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 multdestination traffic. Switches that do not separate unicast and multdestination 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 multdestination 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](#) | 26

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 28 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

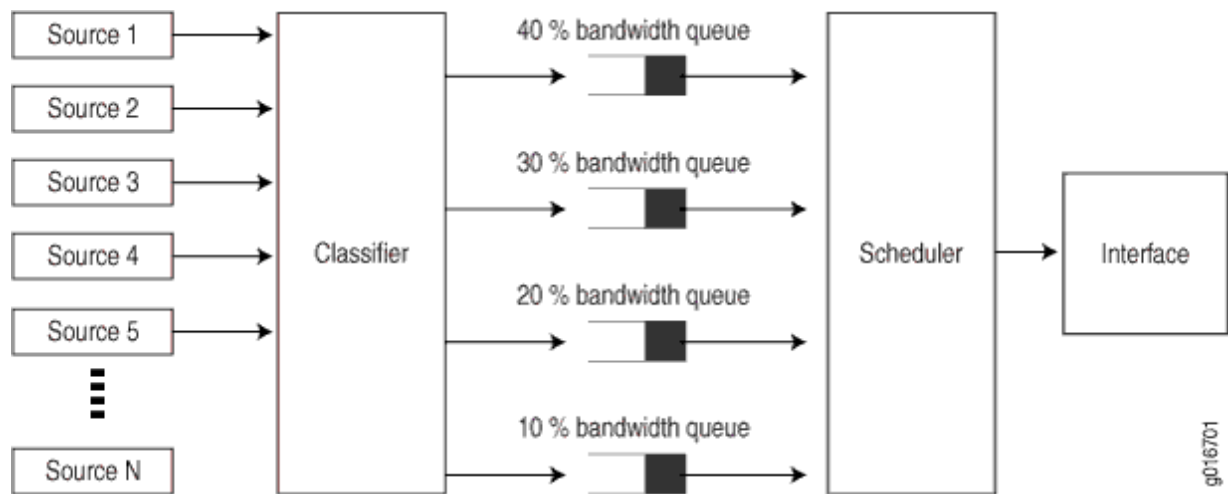
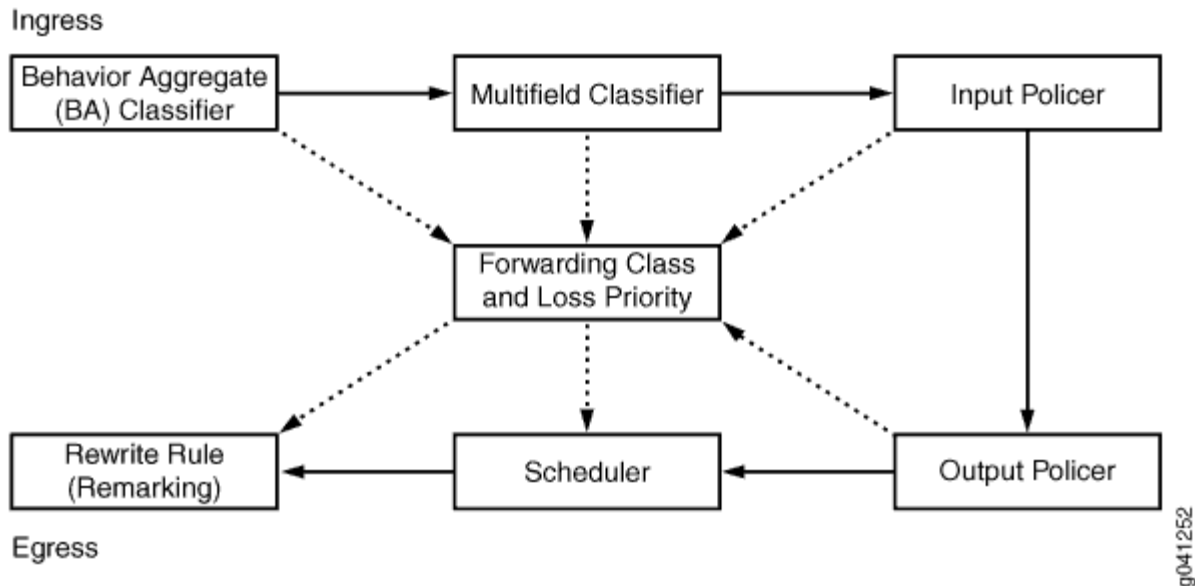


Figure 4 on page 29 shows the packet flow through the CoS components that you can configure.

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 | 30](#)
- [Default Forwarding Class Sets \(Priority Groups\) | 31](#)
- [Default Code-Point Aliases | 32](#)
- [Default Classifiers | 34](#)
- [Default Rewrite Rules | 39](#)
- [Default Drop Profile | 39](#)
- [Default Schedulers | 39](#)
- [Default Scheduler Maps | 43](#)
- [Default Shared Buffer Configuration | 43](#)

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 30](#) 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 32](#) 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 32](#) 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 34](#) 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 35](#) 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 36](#) 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 37](#) 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 38](#) 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 39](#) shows the default drop profile configuration.

Table 11: Default Drop Profile

Fill Level	Drop Probability
100	100

Default Schedulers

[Table 12 on page 40](#) 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 40](#), 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 43 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 43 and Table 15 on page 44 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[Overview of Junos OS CoS | 2](#)[Understanding Junos CoS Components | 21](#)[Understanding Default CoS Scheduling and Classification | 321](#)[Understanding CoS Classifiers | 96](#)[Understanding CoS Classifiers](#)[Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130](#)[Understanding CoS Code-Point Aliases | 90](#)[Understanding CoS Forwarding Classes | 155](#)[Understanding CoS Rewrite Rules | 125](#)[Understanding CoS Output Queue Schedulers | 338](#)[Understanding CoS Port Schedulers | 368](#)[Understanding CoS WRED Drop Profiles | 276](#)

CoS Support on QFX Series Switches, EX4600 Line of Switches, and QFabric Systems

IN THIS SECTION

- [CoS Feature Support | 45](#)
- [Classifier and Rewrite Rule Ethernet Interface Type Support | 49](#)
- [CoS Operational Comparison Between QFX5100, QFX5120, QFX5130, QFX5200, QFX5210, QFX5220, and QFX5700 Switches | 52](#)
- [QFX10000 Switch Classifier and Rewrite Rule Support \(Scaling\) | 57](#)

Juniper Networks data center switches differ in some aspects of class-of-service (CoS) support because of differences in the way the switches are used in networks, and because of hardware differences such as different chipsets or different interface capabilities.

This topic summarizes CoS support on QFX Series switches, the EX4600 line of switches, and QFabric systems.

CoS Feature Support

The first two tables list CoS feature support for newer ELS-CLI-based platforms ([Table 16 on page 45](#)) such as the QFX5000 line, the EX4600 line, and QFX10000 switches, and for legacy-CLI-based platforms ([Table 17 on page 48](#)) such as QFX3500 switches and QFabric systems. Some legacy-CLI-based platforms can also run the ELS CLI.

Table 16: QFX10000, QFX5000 Line, and EX4600 Line CoS Features

Feature	QFX10000	QFX 5000 Line, EX4600 Line	QFX5220/QFX5130/ QFX5700
Class of service (CoS)—Class-based queuing with prioritization	Yes	Yes	Yes
CoS—Separate unicast and multi-destination classifiers, forwarding classes, and output queues	No	Yes	Yes (except multi-destination classifiers. Use firewall filters to classify multicast traffic.)
CoS—Shared unicast and multideestination classifiers, forwarding classes, and output queues	Yes	No	No
CoS support on link aggregation groups (LAGs)	Yes	Yes	Yes

Table 16: QFX10000, QFX5000 Line, and EX4600 Line CoS Features *(Continued)*

Feature	QFX10000	QFX 5000 Line, EX4600 Line	QFX5220/QFX5130/ QFX5700
Enhanced transmission selection (ETS) hierarchical port scheduling	Yes (starting in Junos OS Release 17.3)	QFX5100, QFX 5110, EX4600— Yes QFX5120, QFX5200, QFX5210, EX4650 —No	No
Direct port scheduling	Yes	Yes	Yes
Queue shaping	Yes NOTE: Uses the transmit-rate statement with the exact option.	Yes NOTE: Uses the shaping-rate statement.	Yes
Explicit congestion notification (ECN)	Yes	Yes	Yes
Priority-based flow control (PFC)	Yes	Yes	Yes
Re-marking of bridged packets	Yes	Yes	Yes
Weighted random early detection (WRED) packet drop profiles and tail drop	Yes	Yes	Yes
802.3X Ethernet PAUSE	Yes	Yes	No
Layer 2 ingress packet classification and egress rewrite rules	Yes	Yes	Yes

Table 16: QFX10000, QFX5000 Line, and EX4600 Line CoS Features (*Continued*)

Feature	QFX10000	QFX 5000 Line, EX4600 Line	QFX5220/QFX5130/ QFX5700
MPLS EXP ingress packet classification and egress rewrite rules	Yes	Yes	No
Layer 3 ingress packet classification and egress rewrite rules	Yes	Yes	Yes (Both IPv4 and IPv6 traffic must share the same classifier.)
Virtual output queue (VOQ) architecture	Yes	No	No
Software shared buffer configurability	No (uses VOQ)	Yes	Yes, with the following restrictions: <ul style="list-style-type: none"> multicast partition is not supported in the egress shared buffer pool. See <i>buffer-partition (Egress)</i>. lossy and lossless partitions must have the same percentage values for ingress and egress shared buffer pools.
Shared buffer Alpha configurability	No	Yes	Yes
Buffer monitoring	No	Yes	Yes
CoS command to detect the source of RED-dropped packets	Yes	No	No

Table 17 on page 48 shows CoS support for legacy-CLI-based switches.

Table 17: QFX3500 and QFX3600 Switch, and QFabric System CoS Features (As of Software Release 15.1X53-D30)

Feature	QFX3500	QFX3600	QFabric System
Class of service (CoS)—Class-based queuing with prioritization	Yes	Yes	Yes
CoS—Separate unicast and multideestination classifiers, forwarding classes, and output queues	Yes	Yes	Yes
CoS support on link aggregation groups (LAGs)	Yes	Yes	Yes
Enhanced transmission selection (ETS) hierarchical port scheduling	Yes	Yes	Yes
Direct port scheduling	No	No	No
Queue shaping	Yes	Yes	Yes
Explicit congestion notification (ECN)	Yes	Yes	Yes
Priority-based flow control (PFC)	Yes	Yes	Yes
Re-marking of bridged packets	Yes	Yes	Yes
Priority remapping on native Fibre Channel interfaces	Yes	No	No
Weighted random early detection (WRED) tail-drop profiles	Yes	Yes	Yes
802.3X Ethernet PAUSE	Yes	Yes	Yes
Layer 2 ingress packet classification and egress rewrite rules	Yes	Yes	Yes
MPLS EXP ingress packet classification and egress rewrite rules	Yes	Yes	Yes
Layer 3 ingress packet classification and egress rewrite rules	Yes	Yes	Yes

Table 17: QFX3500 and QFX3600 Switch, and QFabric System CoS Features (As of Software Release 15.1X53-D30) (Continued)

Feature	QFX3500	QFX3600	QFabric System
Software buffer configurability	Yes	Yes	No

Classifier and Rewrite Rule Ethernet Interface Type Support

The next two tables in this topic list CoS Ethernet support for classifiers and rewrite rules on different interface types for QFX10000 switches ([Table 18 on page 49](#)), and for QFX5100, QFX5110, QFX5120, QFX5200, QFX5210, QFX5220, QFX3500, QFX3600, EX4600, and EX4650 switches, and QFabric systems ([Table 19 on page 50](#)).

On 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 18 on page 49](#) shows on which interfaces you can configure and apply classifiers and rewrite rules.

Table 18: 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

Table 18: Ethernet Interface Support for Classifier and Rewrite Rule Configuration (QFX10000 Switches) (Continued)

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

On QFX5100, QFX5110, QFX5120, QFX5200, QFX5210, QFX3500, QFX3600, EX4600, and EX4650 switches, and QFabric systems, you cannot apply classifiers or rewrite rules to Layer 2 physical interfaces or to Layer 3 logical interfaces. [Table 19 on page 50](#) shows on which interfaces you can configure and apply classifiers and rewrite rules.

Table 19: Ethernet Interface Support for Classifier and Rewrite Rule Configuration (QFX5100, QFX5110, QFX5120, QFX5200, QFX5210, EX4600, EX4650, QFX3500, and QFX3600 Switches, and QFabric Systems)

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interface (Unit 0 Only)	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 19: Ethernet Interface Support for Classifier and Rewrite Rule Configuration (QFX5100, QFX5110, QFX5120, QFX5200, QFX5210, EX4600, EX4650, QFX3500, and QFX3600 Switches, and QFabric Systems) (Continued)

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interface (Unit 0 Only)	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 multidestination and DSCP multidestination classifiers are applied to all interfaces and cannot be applied to individual interfaces. No DSCP IPv6 multidestination classifier is supported. IPv6 multidestination traffic uses the DSCP multidestination classifier.

On QFX5220, QFX5130, and QFX5700 switches, you cannot apply classifiers or rewrite rules to Layer 2 or Layer 3 physical interfaces. [Table 20 on page 51](#) shows on which interfaces you can configure and apply classifiers and rewrite rules.

Table 20: Ethernet Interface Support for Classifier and Rewrite Rule Configuration (QFX5220, QFX5130, and QFX5700 Switches)

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interfaces	Layer 3 Physical Interfaces	Layer 3 Logical Interfaces
Fixed classifier	No	Yes	No	Yes
DSCP classifier	No	Yes	No	Yes

Table 20: Ethernet Interface Support for Classifier and Rewrite Rule Configuration (QFX5220, QFX5130, and QFX5700 Switches) (Continued)

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interfaces	Layer 3 Physical Interfaces	Layer 3 Logical Interfaces
DSCP IPv6 classifier	No	No	No	No
IEEE 802.1p classifier	No	Yes	No	Yes
EXP classifier	No	No	No	No
DSCP rewrite rule	No	Yes	No	Yes
DSCP IPv6 rewrite rule	No	No	No	No
IEEE 802.1p rewrite rule	No	Yes	No	Yes
EXP rewrite rule	No	No	No	No

NOTE: QFX5220, QFX5130, and QFX5700 switches do not support DSCP IPV6 classifiers and rewrite rules. Instead, attach DSCP classifiers and rewrite rules on family `inet6`.

CoS Operational Comparison Between QFX5100, QFX5120, QFX5130, QFX5200, QFX5210, QFX5220, and QFX5700 Switches

CoS feature support is mostly the same for QFX5100, QFX5120, QFX5130, QFX5200, QFX5210, QFX 5220, QFX5700 switches, but there are some CoS operational differences due to different chipsets among these platforms. [Table 21 on page 53](#) details both the similarities and differences for CoS on QFX5100, QFX5120, QFX5200, QFX5210, and QFX5220 switches.

Table 21: CoS Operational Comparison Between QFX5100, QFX5120, QFX5130, QFX5200, QFX5210, QFX5220, and QFX5700 Switches (Continued)

CoS Feature	QFX5100	QFX5120	QFX5130/ QFX5700	QFX5200	QFX5210	QFX5220	Change in Operation
Queuing and Scheduling	LLS and three- level hierarchy	Fixed hierarchical scheduling (FHS) and two-level hierarchy	Fixed hierarchical scheduling (FHS) and two-level hierarchy	Fixed hierarchical scheduling (FHS) and two- level hierarchy	Fixed hierarchical scheduling (FHS) and two-level hierarchy	Fixed hierarchical scheduling (FHS) and two-level hierarchy	ETS and FC- Set are not supported on QFX5120, QFX5130, QFX5200, QFX5210, QFX5220, and QFX5700 due to FHS.
# Unicast Queues	8	8	8	8	8	8	N/A
# Multicast Queues	4	2	4	2	2	2	N/A
CPU Queues	44	44	44	44	44	44	N/A
Host Path Scheduling	48 queues directly attached to port	48 queues attached to L0	48 queues attached to L0	48 queues attached to L0	48 queues attached to L0	48 queues attached to L0	No customer visible change.
FC2Q	4 profiles	4 profiles	4 profiles	4 profiles	4 profiles	4 profiles	N/A
DSCP classifier table	128 profiles	128 profiles	64 profiles	128 profiles	128 profiles	64 profiles	N/A

Table 21: CoS Operational Comparison Between QFX5100, QFX5120, QFX5130, QFX5200, QFX5210, QFX5220, and QFX5700 Switches (Continued)

[illegible]

Table 21: CoS Operational Comparison Between QFX5100, QFX5120, QFX5130, QFX5200, QFX5210, QFX5220, and QFX5700 Switches (Continued)

CoS Feature	QFX5100	QFX5120	QFX5130/ QFX5700	QFX5200	QFX5210	QFX5220	Change in Operation
Queueing Levels	Four levels physical queue level, logical queue level, CoS level, and port level	Three levels, logical queue level, CoS level, and port level.	Three levels, logical queue level, CoS level, and port level.	Three levels, logical queue level, CoS level, and port level.	Three levels, logical queue level, CoS level, and port level.	Three levels, logical queue level, CoS level, and port level.	N/A
Multidestination Traffic	Default scheduler map reserves 20% bandwidth for multicast and 80% of unicast traffic reserved between BE, FCoE, NoLoss and NC traffic types.	Same as QFX5100 switches	By default all multicast traffic mapped to Q8. Q8 is given 20% bandwidth in default scheduler. To classify multicast traffic to different queues (Q9,10,11) use firewall filters.	Each level 0 node is receiving both multicast and unicast traffic, so it is not possible to differentiate at the port level to apply shaping on multicast traffic.	Same as QFX5200 switches	By default all multicast traffic mapped to Q8. Q8 is given 20% bandwidth in default scheduler. To classify multicast traffic to different queue (Q9) use firewall filters.	N/A

The following limitations on QFX5200 and QFX5210 switches do not exist on QFX5100 switches.

- CoS flexible hierarchical scheduling (ETS) is not supported on QFX5200 or QFX5210 switches.
- QFX5200 and QFX5210 switches support only one queue with strict-high priority because these switches do not support flexible hierarchical scheduling.

NOTE: QFX5100 switches support multiple queues with strict-high priority when you configure a forwarding class set.

- QFX5200 CoS policers do not support global management counters accessed by all ports. Only management counters local to a pipeline are supported—this means that QFX5200 management counters work only on traffic received on ports that belong to the pipeline in which the counter is created.
- Due to the cross-point architecture on QFX5200 and QFX5210 switches, all buffer usage counters are maintained separately. When usage counters are displayed with the command `show class-of-service shared-buffer`, various pipe counters are displayed separately.
- On QFX5200 and QFX5210 switches, port schedulers are supported instead of FC-SET.
- On QFX5200 and QFX5210 switches, it is not possible to group multiple forwarding classes into a forwarding class set (fc-set) and apply output traffic control profile on the fc-set. ETS for an fc-set is not supported. Because each L0 node schedules both the unicast and multicast queue of L1 node, it is not possible to differentiate multicast and unicast traffic at the port level and apply minimum bandwidth between unicast and multicast. It can only be supported at CoS level L0.
- Because QFX5200 and QFX5210 switches do not support flexible hierarchical scheduling, it is not possible to apply a traffic control profile for a group of forwarding classes.

QFX10000 Switch Classifier and Rewrite Rule Support (Scaling)

You can configure enough classifiers on QFX10000 switches to handle most, if not all, network scenarios. [Table 22 on page 57](#) shows how many of each type of classifiers you can configure, and how many entries you can configure per classifier.

Table 22: Classifier Support by Classifier Type on QFX10000 Switches

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

Table 22: Classifier Support by Classifier Type on QFX10000 Switches (Continued)

Classifier Type	Default Classifier Name	Maximum Number of Classifiers	Maximum Number of Entries per Classifier
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).

There are no default rewrite rules. You can configure enough rewrite rules on QFX10000 switches to handle most, if not all, network scenarios. [Table 23 on page 58](#) shows how many of each type of rewrite rule you can configure, and how many entries you can configure per rewrite rule.

Table 23: Rewrite Rule Support by Rewrite Rule Type on QFX10000 Switches

Rewrite Rule Type	Maximum Number of Rewrite Rule Sets	Maximum Number of Entries per Rewrite Rule Set
IEEE 802.1p (Layer 2)	64	128
DSCP (Layer 3)	32	128
DSCP IPv6 (Layer 3)	32	128
EXP (MPLS)	64	128

RELATED DOCUMENTATION

[Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130](#)

[Understanding Default CoS Scheduling and Classification | 321](#)

[Understanding CoS Classifiers | 96](#)

[Understanding CoS Output Queue Schedulers | 338](#)

[Understanding CoS Virtual Output Queues \(VOQs\) | 406](#)

[Understanding CoS Port Schedulers | 368](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

[Understanding CoS WRED Drop Profiles | 276](#)

CoS on Interfaces

IN THIS CHAPTER

- CoS Inputs and Outputs Overview | 60
- CoS on Virtual Chassis Switch Ports | 61
- CoS on Virtual Chassis Fabric (VCF) EX4300 Leaf Devices (Mixed Mode) | 66
- Understanding CoS on OVSDB-Managed VXLAN Interfaces | 73
- Configuring CoS on OVSDB-Managed VXLAN Interfaces | 78
- Assigning CoS Components to Interfaces | 87

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 24 on page 60](#).

Table 24: CoS Mappings—Inputs and Outputs

CoS Mappings	Inputs	Outputs	Comments
classifiers	code-points	forwarding-class, loss-priority	The map sets the forwarding class and packet loss priority (PLP) for a specific set of code points.
drop-profile-map	loss-priority, protocol	drop-profile	The map sets the drop profile for a specific PLP and protocol type.
rewrite-rules	loss-priority, forwarding-class	code-points	The map sets the code points for a specific forwarding class and PLP.

Table 24: CoS Mappings—Inputs and Outputs *(Continued)*

CoS Mappings	Inputs	Outputs	Comments
rewrite-value (Fibre Channel Interfaces)	<i>forwarding-class</i>	<i>code-point</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](#) | 26

CoS on Virtual Chassis Switch Ports

IN THIS SECTION

- [Access Interface CoS Support](#) | 62
- [VCP Interface CoS Support](#) | 63
- [CPU-Generated Host Outbound Traffic](#) | 65

QFX Series and EX4600 Virtual Chassis devices have access ports to connect to external peer devices. Virtual Chassis devices also have Virtual Chassis ports (VCPs) to interconnect members of the Virtual Chassis, in a similar way that QFabric system Node devices have fabric (fte) ports to connect to the QFabric system Interconnect device. VCPs are not used for external access.

Class of service (CoS) on Virtual Chassis access ports is the same as CoS on these devices when they are in standalone mode or used as QFabric system Node devices. However, CoS on VCPs differs in several ways from CoS on QFabric system Node device fabric ports.

This topic describes CoS support on Virtual Chassis access interfaces and on VCPs.

Access Interface CoS Support

CoS on Virtual Chassis access interfaces is the same as CoS on standalone device and Node device access interfaces, except for shared buffer settings. The documentation for QFX Series and EX4600 switch CoS on access interfaces applies to Virtual Chassis access interfaces, except some of the shared buffer documentation.

Similarities in CoS Support on Virtual Chassis Access Interfaces Compared to Standalone Device (or QFabric system Node device) Access Interfaces

Virtual Chassis access interfaces support the following CoS features in the same way as access interfaces on standalone devices and QFabric system Node devices:

- Forwarding classes—The default forwarding classes, queue mapping, and packet drop attributes ([Table 25 on page 62](#)) are the same:

Table 25: Default Forwarding Class Configuration

Default Forwarding Class	Default Queue Mapping	Default Packet Drop Attribute
best-effort (be)	0	drop
fcoe	3	no-loss
no-loss	4	no-loss
network-control (nc)	7	drop
mcast	8	drop

- Packet classification—Classifier default settings and configuration are the same. Support for behavior aggregate, multifield, multideestination, and fixed classifiers is the same.
- Enhanced transmission selection (ETS)—This data center bridging (DCB) feature that supports hierarchical scheduling has the same defaults and user configuration, including forwarding class set (priority group) and traffic control profile configuration.
- Priority-based flow control (PFC)—This DCB feature that supports lossless transport has the same defaults and user configuration, including support for six lossless priorities (forwarding classes).
- Ethernet PAUSE—This feature has the same defaults and configuration.

- Queue scheduling—This feature has the same defaults, configuration, and scheduler-to-forwarding-class mapping. Queue scheduling is a subset of hierarchical scheduling.
- Priority group (forwarding class set) scheduling—This feature has the same defaults and configuration. Priority group scheduling is a subset of hierarchical scheduling.
- WRED profiles—This feature has the same defaults and configuration.
- Code-point aliases—This feature has the same defaults and configuration.
- Rewrite rules—This feature has the same defaults and configuration (no default rewrite rules applied to egress traffic).
- Host outbound traffic—This feature has the same defaults and configuration.

Differences in CoS Support on Virtual Chassis Access Interfaces Compared to Standalone Device (or QFabric system Node device) Access Interfaces

The default shared buffer settings and the way in which you configure shared buffers are the same on Virtual Chassis access interfaces as on standalone and QFabric system Node devices. The difference is that on Virtual Chassis access interfaces, the shared buffer configuration is global and applies to all access ports on all members of the Virtual Chassis, while on standalone or QFabric system Node devices, you can configure different buffer settings on different access interfaces.

You cannot configure different shared buffer settings for different Virtual Chassis members. All members of a Virtual Chassis use the same shared buffer configuration.

VCP Interface CoS Support

CoS on the VCP interfaces that connect the Virtual Chassis members is similar to CoS on the fabric interfaces of QFabric system Node devices, but there are several important differences:

Similarities in CoS Support on VCP Interfaces and QFabric System Node Device Fabric Interfaces

VCP interfaces support full hierarchical scheduling (ETS). ETS includes the following CoS features. VCP interfaces support no other CoS features.

- Creating forwarding class sets (priority groups) and mapping forwarding classes to forwarding class sets.
- Scheduling individual output queues. The scheduler defaults and configuration are the same as the scheduler on access interfaces.

- Scheduling priority groups (forwarding class sets) using a traffic control profile. The defaults and configuration are the same as on access interfaces.

NOTE: You cannot attach classifiers, congestion notification profiles, or rewrite rules to VCP interfaces. You cannot attach scheduler maps to VCP interfaces on QFX platforms that do not support ETS. Also, you cannot configure buffer settings on VCP interfaces. You can only attach forwarding class sets and traffic control profiles to VCP interfaces (as well as scheduler maps *if* the platform supports ETS).

The behavior of lossless traffic across 40-Gigabit VCP interfaces is the same as the behavior of lossless traffic across QFabric system Node device fabric ports. The system automatically enables flow control for lossless forwarding classes (priorities). The system dynamically calculates buffer headroom that is allocated from the global lossless-headroom buffer for the lossless forwarding classes on each 40-Gigabit VCP interface. If there is not enough global lossless-headroom buffer space to support the number of lossless flows on a 40-Gigabit VCP interface, the system generates a syslog message.

NOTE: After you configure lossless transport on a Virtual Chassis, check the syslog messages to ensure that there is sufficient buffer space to support the configuration.

NOTE: If you break out a 40-Gigabit VCP interface into 10-Gigabit VCP interfaces, lossless transport is not supported on the 10-Gigabit VCP interfaces. Lossless transport is supported only on 40-Gigabit VCP interfaces. (10-Gigabit access interfaces support lossless transport.)

Differences in CoS Support on VCP Interfaces and QFabric System Node Device Fabric Interfaces

Although most of the CoS behavior on VCP interfaces is similar to CoS behavior on the fabric ports of QFabric system Node devices, there are some important differences:

- Hierarchical scheduling (queue and priority group scheduling)—On QFabric system Node device fabric interfaces, you can apply a different hierarchical scheduler (traffic control profile) to different priority groups (forwarding class sets) on different interfaces. However, on VCP interfaces, the schedulers that you apply to priority groups are global to all VCP interfaces. One hierarchical scheduler controls scheduling for a priority group on all VCP interfaces.

You attach a scheduler to VCP interfaces using the global identifier (vcp-*) for VCP interfaces. For example, if you want to apply a traffic control profile (traffic control profiles contain both queue and

priority group scheduling configuration) named *vcp-hpc-tcp* to a forwarding class set named *vcp-hpc-fcset*, you include the following statement in the configuration:

```
[edit]
user@switch# set class-of-service interfaces vcp-* forwarding-class-set vcp-hpc-fcset output-traffic-control-profile vcp-hpc-tcp
```

The system applies the hierarchical scheduler *vcp-hpc-tcp* to the traffic mapped to the priority group *vcp-hpc-fcset* on all VCP interfaces.

- You cannot attach classifiers, congestion notification profiles, or rewrite rules to VCP interfaces. Also, you cannot configure buffer settings on VCP interfaces. Similar to QFabric system Node device fabric interfaces, you can only attach forwarding class sets and traffic control profiles to VCP interfaces.
- Lossless transport is supported only on 40-Gigabit VCP interfaces. If you break out a 40-Gigabit VCP interface into 10-Gigabit VCP interfaces, lossless transport is not supported on the 10-Gigabit VCP interfaces.

CPU-Generated Host Outbound Traffic

CPU-generated host outbound traffic is forwarded on the network-control forwarding class, which is mapped to queue 7. If you use the default scheduler, the network-control queue receives a guaranteed minimum bandwidth (transmit rate) of 5 percent of port bandwidth. The guaranteed minimum bandwidth is more than sufficient to ensure lossless transport of host outbound traffic.

However, if you configure and apply a scheduler instead of using the default scheduler, you must ensure that the network-control forwarding class (or whatever forwarding class you configure for host outbound traffic) receives sufficient guaranteed bandwidth to prevent packet loss.

TIP: If you configure a scheduler instead of using the default scheduler, we recommend that you configure the network-control queue (or the queue you configure for host outbound traffic if it is not the network-control queue) as a strict-high priority queue. Strict-high priority queues receive the bandwidth required to transmit their entire queues before other queues are served. To limit the amount of bandwidth a strict-high priority queue can consume (and to prevent the strict-high priority queue from starving other queues), apply a shaping rate to the strict-high priority traffic in the scheduler configuration.

As with all strict-high priority traffic, if you configure the network-control queue (or any other queue) as a strict-high priority queue, you must also create a separate forwarding class set (priority group) that contains only strict-high priority traffic, and apply the strict-high priority forwarding class set and its traffic control profile (hierarchical scheduler) to the VCP interfaces.

RELATED DOCUMENTATION

Understanding Default CoS Settings	 30
Understanding CoS Classifiers	 96
Understanding CoS Forwarding Classes	 155
Understanding CoS Hierarchical Port Scheduling (ETS)	 438
Understanding CoS Buffer Configuration	 684
Understanding CoS WRED Drop Profiles	 276
Understanding CoS Rewrite Rules	 125
Understanding CoS Flow Control (Ethernet PAUSE and PFC)	 220
Understanding Host Routing Engine Outbound Traffic Queues and Defaults	 271

CoS on Virtual Chassis Fabric (VCF) EX4300 Leaf Devices (Mixed Mode)

IN THIS SECTION

- [VCF CoS in Mixed Mode with an EX4300 Leaf Device](#) | 67
- [Scheduling on an EX4300 VCF Leaf Device](#) | 69

A Virtual Chassis Fabric (VCF) uses QFX5100 switches as spine devices and can use QFX5100, QFX3500, QFX3600, and EX4300 switches as leaf devices. When a VCF includes more than one type of leaf device (mixed mode), the CoS feature support on the VCF depends on the capability of the lowest-featured device. In mixed mode, the supported CoS features are the “lowest common denominator” of the features supported by the leaf devices. If one leaf device does not support a particular feature, that feature is not supported on the VCF even if every other leaf device supports the feature.

NOTE: EX4300 leaf devices do not support several CoS features that are supported on QFX5100, QFX3600, and QFX3500 devices. However, even when a VCF includes an EX4300 leaf device, other leaf devices might support those CoS features.

VCF CoS in Mixed Mode with an EX4300 Leaf Device

In mixed mode, if all of the leaf devices are QFX5100, QFX3500, and QFX3600 switches, the full QFX Series CoS feature set is available, including data center bridging (DCB) features such as enhanced transmission selection (ETS, IEEE 802.1Qaz), priority-based flow control (PFC, IEEE 802.1Qbb), and Data Center Bridging Exchange Protocol (DCBX, an extension of LLDP, IEEE 802.1AB).

However, the EX4300 leaf device does not support DCB standards (ETS, PFC, DCBX). The lack of support for DCB standards means that the EX4300 leaf device does not support lossless transport. So a VCF that includes an EX4300 as a leaf device does not support lossless storage traffic such as Fibre Channel over Ethernet (FCoE).

In addition, a VCF with an EX4300 leaf device either does not support or has limited support for some other CoS features that the QFX Series switches support, including some buffer configuration features, some packet rewrite features, and Ethernet PAUSE (IEEE 802.3X).

[Table 26 on page 67](#) summarizes the CoS support on a VCF in mixed mode with one or more EX4300 leaf devices.

Table 26: Support of QFX CoS Features on a VCF in Mixed Mode with an EX4300 Leaf Device

QFX Series CoS Feature	Support in Mixed Mode with an EX4300 Leaf Device
Forwarding Classes	The EX4300 leaf device uses the QFX Series default forwarding classes, the default QFX Series forwarding class to queue mapping, and the QFX Series maximum number of supported forwarding classes (12).
Lossless Forwarding Classes	Not supported. For example, the QFX Series default lossless forwarding classes fcoe and no-loss are not treated as lossless forwarding classes. Traffic mapped to lossless forwarding classes (default lossless forwarding classes or user-defined lossless forwarding classes) is treated as best-effort traffic.
Shared buffer configuration	Ingress shared buffer configuration is not supported. Egress shared buffer configuration does not support partitioning into three buffer pools. If there is a shared buffer configuration, only the total egress shared buffer configuration is used. Ingress shared buffer configuration and egress buffer partitioning configuration is ignored.
Classifier on a Layer 2 interface	One classifier per protocol is supported on a port. On a physical port, for a particular protocol, the same Layer 2 classifier is used on all of the logical interfaces.

Table 26: Support of QFX CoS Features on a VCF in Mixed Mode with an EX4300 Leaf Device
(Continued)

QFX Series CoS Feature	Support in Mixed Mode with an EX4300 Leaf Device
Classifier on a Layer 3 interface	Supported.
Multi-destination classifier	<p>Supported.</p> <p>The EX4300 leaf device uses the same default classifier as the QFX5100 spine device. As on QFX Series switches, a multi-destination classifier is global and is applied to all VCF interfaces. Multi-destination classifiers are valid only for multicast forwarding classes. You can configure two multi-destination classifiers, one for IEEE 802.1p traffic and one for DSCP traffic (the DSCP multi-destination classifier applies to both IPv4 and IPv6 traffic).</p>
Congestion notification profile	<p>Not supported.</p> <p>If a congestion notification profile is configured on the QFX5100 spine device, it is ignored because the EX4300 leaf device does not support lossless transport, so end-to-end lossless behavior is not possible</p>
Ethernet PAUSE (IEEE 802.3X)	<p>Not supported.</p> <p>If Ethernet PAUSE is configured, it is ignored.</p>
Hierarchical scheduling (ETS)	<p>Translated into port-based scheduling.</p> <p>The EX4300 device does not support ETS scheduling. A VCF translates ETS scheduling configured on a QFX5100 spine device into port scheduling on an EX4300 leaf device. The hierarchical structure of mapping forwarding classes into forwarding class sets (fc-sets) is ignored.</p> <p>"Scheduling on an EX4300 VCF Leaf Device" on page 69 provides details on how a VCF translates QFX Series ETS scheduling into port scheduling on an EX4300 leaf device.</p>
Hierarchical scheduling (ETS) on a spine device VCP port	On QFX5100 VCP ports, the hierarchical mapping of forwarding classes to forwarding class sets is supported. However, scheduling on an EX4300 leaf device is translated into port scheduling.

Table 26: Support of QFX CoS Features on a VCF in Mixed Mode with an EX4300 Leaf Device
(Continued)

QFX Series CoS Feature	Support in Mixed Mode with an EX4300 Leaf Device
Drop profile (WRED)	<p>QFX Series drop profiles are supported. The EX4300 device as a standalone switch supports four packet loss priorities. However, as part of a mixed mode VCF, the EX4300 leaf device supports only the three packet loss priorities that the QFX Series switches support:</p> <ul style="list-style-type: none"> • low • medium-high • high <p>Supporting only three packet loss priorities means that the behavior of the EX4300 switch as a leaf device is different from the behavior as a standalone switch.</p>
Rewrite rules on a Layer 2 interface	Supported, but with a limit of one rewrite rule per physical interface. All traffic uses the same rewrite rule.
Rewrite rules on a Layer 3 interface	Supported, but with a limit of one rewrite rule per physical interface. The same rewrite rule is used on all traffic on the interface.
Rewrite value for FCoE traffic	<p>Not supported.</p> <p>If a rewrite value for FCoE traffic, is configured, it is ignored. (A mixed mode VCF does not support lossless traffic.)</p>

In addition to the CoS limitations shown in [Table 26 on page 67](#), using wild cards in a LAG configuration is not supported in mixed mode with one or more EX4300 leaf devices.

Scheduling on an EX4300 VCF Leaf Device

Because the EX4300 leaf device does not support ETS, the VCF translates the ETS scheduling configuration into the port scheduling configuration that the EX4300 device supports. The QFX5100 spine device uses two-tier ETS scheduling, as described in detail in ["Understanding CoS Hierarchical Port Scheduling \(ETS\)" on page 438](#).

Briefly, ETS allocates port bandwidth into forwarding class sets (priority groups) and forwarding classes (priorities) in a hierarchical manner. Each forwarding class set consists of individual forwarding classes, with each forwarding class mapped to an output queue.

Port bandwidth (minimum guaranteed bandwidth and maximum bandwidth) is allocated to each forwarding class set. Forwarding class set bandwidth is in turn allocated to the forwarding classes in the forwarding class set. If a forwarding class does not use its bandwidth allocation, other forwarding classes within the same forwarding class set can share the unused bandwidth. If the forwarding classes in a forwarding class set do not use the bandwidth allocated to that forwarding class set, other forwarding class sets on the port can share the unused bandwidth. (This is how ETS increases port bandwidth utilization, by sharing unused bandwidth among forwarding classes and forwarding class sets.)

However, the EX4300 leaf device supports port scheduling, not ETS. Port scheduling is a “flat” scheduling method that allocates bandwidth directly to forwarding classes in a non-hierarchical manner.

The VCF translates the two tiers of the ETS scheduling configuration (forwarding class sets and forwarding classes) into a single port scheduling configuration as follows:

- The bandwidth allocated to a forwarding class set is divided equally among the forwarding classes in the forwarding class set. (Traffic control profiles schedule bandwidth allocation to forwarding class sets.) The minimum guaranteed bandwidth (guaranteed-rate) and maximum bandwidth limit (shaping-rate) of the forwarding class set determine the guaranteed minimum bandwidth and the maximum bandwidth the forwarding classes receive, *unless* those values are different in the forwarding class scheduler configuration.
- If there is an explicit forwarding class bandwidth scheduler configuration, it overrides the forwarding class set configuration. Bandwidth scheduling values that are not explicitly configured in a forwarding class scheduler use the values from the forwarding class set (the traffic control profile configuration). Forwarding class schedulers control the minimum guaranteed bandwidth (transmit-rate), the maximum bandwidth (shaping-rate), and the priority (priority) for each forwarding class (output queue). Because the priority value is not configured at the forwarding class set level, the priority configured in the forwarding class scheduler is always used.

The following two scenarios illustrate how a VCF translates an ETS configuration into a port scheduling configuration:

Scenario 1

A forwarding class set named `fc-set-1` has a configured guaranteed minimum bandwidth (guaranteed-rate) of 4G, and a configured maximum bandwidth (shaping-rate) of 5G.

Forwarding class set `fc-set-1` consists of two forwarding classes, named `fc-1` and `fc-2`:

- Forwarding class `fc-1` has a guaranteed minimum bandwidth (transmit-rate) of 2.5G. There is no configured maximum bandwidth (shaping-rate).
- Forwarding class `fc-2` has a guaranteed minimum bandwidth (transmit-rate) of 1.5G. There is no configured maximum bandwidth (shaping-rate).

On the EX4300 leaf device, the ETS configuration above is translated approximately to the following port scheduling configuration:

- **Guaranteed minimum bandwidth**—Because guaranteed minimum bandwidth has been explicitly configured in the forwarding class scheduler, forwarding class fc-1 receives a transmit rate of 2.5G and forwarding class fc-2 receives a transmit rate of 1.5G.

NOTE: If there had been no forwarding class scheduler transmit-rate configuration, then the forwarding class set minimum guaranteed bandwidth of 4G would have been split evenly between the forwarding classes, with each forwarding class receiving a minimum guaranteed bandwidth rate of 2G.

- **Maximum bandwidth**—Because there is no explicit maximum bandwidth (shaping-rate configuration) for the forwarding classes, the forwarding classes that belong to the forwarding class set receive an equal share of the maximum bandwidth configured at the forwarding class set level in the traffic control profile. Because the forwarding class set maximum bandwidth is 5G, forwarding classes fc-1 and fc-2 each receive a maximum bandwidth of 2.5G.

In this scenario, the minimum guaranteed bandwidth and the maximum bandwidth configured at the forwarding class set hierarchy level are achieved on the forwarding classes that belong to the forwarding class set. (This does not always happen, as Scenario 2 shows.) However, unused bandwidth is not shared the same way. For example, if forwarding class fc-1 experienced a burst of traffic at 3.5G, it would be limited to a maximum of 2.5G and traffic would be dropped. Using ETS, if forwarding class fc-2 was not using its allocated maximum bandwidth, then fc-1 could use (share) that unused bandwidth. But flat port scheduling does not share the unused bandwidth.

Scenario 2

A forwarding class set named fc-set-2 has a configured guaranteed minimum bandwidth (guaranteed-rate) of 6G, and a configured maximum bandwidth (shaping-rate) of 9G.

Forwarding class set fc-set-2 consists of three forwarding classes, named fc-3, fc-4, and fc-5:

- Forwarding class fc-3 has a guaranteed minimum bandwidth (transmit-rate) of 1G. There is no configured maximum bandwidth (shaping-rate).
- Forwarding class fc-4 has a maximum bandwidth (shaping-rate) of 2G. There is no configured guaranteed minimum bandwidth (transmit-rate).
- Forwarding class fc-5 has a guaranteed minimum bandwidth (transmit-rate) of 3G. There is no configured maximum bandwidth (shaping-rate).

On the EX4300 leaf device, the ETS configuration above is translated approximately to the following port scheduling configuration:

- **Guaranteed minimum bandwidth**—Two forwarding classes (fc-3 and fc-5) have an explicitly configured transmit rate, and one forwarding class (fc-4) does not. Forwarding classes fc-3 and fc-5 receive the

minimum guaranteed bandwidth configured in their schedulers, so forwarding class fc-3 receives 1G guaranteed minimum bandwidth and forwarding class fc-5 receives 3G guaranteed minimum bandwidth.

Forwarding class fc-4 does not have an explicitly configured transmit rate, so the port derives the minimum guaranteed bandwidth from the forwarding class set guaranteed rate. Forwarding class set fc-set-2 has a minimum guaranteed bandwidth (guaranteed-rate) of 6G, and there are three forwarding classes in the forwarding class set. Forwarding class fc-4 receives an equal share (one third) of the forwarding class set minimum guaranteed bandwidth. So forwarding class fc-4 is allocated a guaranteed minimum bandwidth (transmit-rate) of 2G (6G divided by 3 forwarding classes = 2G).

- **Maximum bandwidth**—Forwarding class fc-4 has an explicitly configured shaping rate, and forwarding classes fc-3 and fc-5 do not. Forwarding class fc-4 receives the maximum bandwidth configured in its scheduler, so forwarding class fc-4 receives a maximum bandwidth of 2G.

Forwarding classes fc-3 and fc-5 do not have explicitly configured shaping rates, so the port derives the maximum bandwidth from the forwarding class set shaping rate. Forwarding class set fc-set-2 has a maximum bandwidth (shaping-rate) of 9G, and there are three forwarding classes in the forwarding class set. Forwarding classes fc-3 and fc-5 each receive an equal share (one third) of the forwarding class set shaping rate. So forwarding classes fc-3 and fc-5 are allocated a maximum bandwidth of 3G each (9G divided by 3 forwarding classes = 3G).

Forwarding class fc-4 receives less maximum bandwidth than forwarding classes fc-3 and fc-5 because the explicitly configured shaping rate for forwarding class fc-4 is only 2G, and the explicit forwarding class configuration overrides the forwarding class set configuration.

NOTE: Scenario 2 shows that in some cases, the guaranteed minimum bandwidth (guaranteed-rate) and the maximum bandwidth (shaping-rate) configured for a forwarding class set might not be achieved at the forwarding class (queue) level. In Scenario 2, forwarding class set fc-set-2 has a shaping rate of 9G, but the sum of the implemented forwarding class shaping rates is only 8G [(3G for fc-3) + (2G for fc-4) + (3G for fc-5)].

RELATED DOCUMENTATION

[Understanding Default CoS Settings | 30](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

Understanding CoS on OVSDB-Managed VXLAN Interfaces

IN THIS SECTION

- [Classifier and Rewrite Rule Interface Support | 74](#)
- [Classifiers on OVSDB-Managed VXLAN Interfaces | 75](#)
- [Rewrite Rules on OVSDB-Managed VXLAN Interfaces | 76](#)
- [Schedulers on OVSDB-Managed VXLAN Interfaces | 77](#)

You can configure class of service (CoS) features on OVSDB-managed VXLAN interfaces on QFX5100 and QFX10000 Series switches. An OVSDB-managed VXLAN interface uses an OVSDB controller to create and manage the VXLAN interfaces and tunnels. OVSDB-managed VXLAN interfaces support:

- Packet classifiers on ingress interfaces. On network-facing interfaces (interfaces that connect to the network, for example, switch interfaces that connect to a VXLAN gateway), you can configure DSCP classifiers. Fixed classifiers, 802.1p classifiers, and MPLS EXP classifiers are not supported on VXLAN interfaces.

NOTE: MF Filters on access-facing interfaces are applied as a group config and not as a normal filter.

- Packet rewrite rules (to change the code point bits of outgoing packets). On network-facing interfaces, you can configure DSCP rewrite rules. Rewrite rules are not supported on access-facing interfaces, and are not supported for IEEE 802.1p code points.

NOTE: Rewrite rules rewrite the DSCP code point on the VXLAN header only. Rewrite rules do not rewrite the DSCP code point on the inner packet header.

- Packet schedulers on egress interfaces. You can configure schedulers on network-facing and access-facing interfaces.

NOTE: You cannot configure CoS on manually configured VXLAN interfaces.

CoS configuration on OVSDB-managed VXLAN interfaces uses the same CLI statements and configuration constructs as CoS configuration on regular Ethernet interfaces. However, feature support differs on OVSDB-managed VXLAN interfaces and regular Ethernet interfaces. The following sections describe the differences between CoS support on OVSDB-managed VXLAN interfaces and regular Ethernet interfaces:

Classifier and Rewrite Rule Interface Support

The switch 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)

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

- Layer 2 physical interfaces. All underlying logical Layer 2 interfaces on the physical interface use the classifier and rewrite rule configuration on the physical interface. All OVSDB-managed VXLAN traffic on the interface uses the same Layer 2 CoS classifiers and rewrite rules.
- Layer 3 physical interfaces if at least one logical Layer 3 interface is configured on the physical interface. All underlying logical Layer 3 interfaces on the physical interface use the classifier and rewrite rule configuration on the physical interface. All OVSDB-managed VXLAN traffic on the interface uses the same Layer 3 CoS classifiers and rewrite rules.

[Table 27 on page 74](#) shows on which interfaces you can configure and apply classifiers and rewrite rules on *network-facing* interfaces.

Table 27: OSVDB-Managed VXLAN Interface Support for Classifier and Rewrite Rule Configuration on Network-Facing Interfaces

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interfaces	Layer 3 Physical Interfaces (If at Least One Logical Layer 3 Interface Is Defined)	Layer 3 Logical Interfaces
Fixed classifier	Not Supported			
DSCP classifier	Yes	No	Yes	No

Table 27: OSVDB-Managed VXLAN Interface Support for Classifier and Rewrite Rule Configuration on Network-Facing Interfaces (Continued)

CoS Classifiers and Rewrite Rules	Layer 2 Physical Interfaces	Layer 2 Logical Interfaces	Layer 3 Physical Interfaces (If at Least One Logical Layer 3 Interface Is Defined)	Layer 3 Logical Interfaces
DSCP IPv6 classifier	Yes	No	Yes	No
IEEE 802.1p classifier	Not Supported			
EXP classifier	Not Supported			
DSCP rewrite rule	Yes	No	Yes	No
DSCP IPv6 rewrite rule	Yes	No	Yes	No
IEEE 802.1p rewrite rule	Not Supported			
EXP rewrite rule	Not Supported			

NOTE: The switch encapsulates packets in VXLAN after packet classification, and before packet rewrite and scheduling.

Classifiers on OVSDb-Managed VXLAN Interfaces

Classifiers map incoming packets to a CoS service level, based on the code points in the header of the incoming packet. At the ingress interface, the switch reads the code point value in the packet header, then assigns the packet to the forwarding class and loss priority mapped to that code point value. The forwarding class is mapped to an egress queue and to scheduling properties. OVSDb-managed VXLAN interfaces support packet classification based on DSCP code points on all ingress interfaces, and packet classification based on DSCP multi-field (MF DSCP) code points or for behavior aggregate (BA) classification, the DSCP, DSCP IPv6, or IP precedence bits of the IP header convey the behavior aggregate class information on access-facing interfaces.

If you do not configure classifiers, the switch uses the default CoS settings to classify incoming traffic, as described in ["Understanding Default CoS Scheduling and Classification" on page 321](#).

Classifier configuration on an OVSDB-managed VXLAN switch interface is similar to classifier configuration on any other type of ingress interface (see ["Understanding CoS Classifiers" on page 96](#)). However, on OVSDB-managed VXLAN interfaces, there is a difference in the way you can apply classifiers to Layer 2 interfaces compared to non-VXLAN interfaces. On OVSDB-managed VXLAN interfaces, you apply classifiers to Layer 2 physical interfaces, and the underlying logical interfaces use the classifier configuration applied on the physical interface. On non-VXLAN interfaces, you apply Layer 2 classifiers to logical interface unit 0 (all other logical interfaces on the port use the classifier configured on unit 0), and not to physical interfaces.

Classifiers on Access-Facing Interfaces

When a packet enters an ingress switch from a server (or other source), you can map it to a forwarding class and a loss priority based on its DSCP multi-field (MF DSCP) classifiers code points. The forwarding class is mapped to an egress queue and to scheduling properties. For behavior aggregate (BA) classification, the DSCP, DSCP IPv6, or IP precedence bits of the IP header convey the behavior aggregate class information.

Classifiers on Network-Facing Interfaces

When a packet enters an egress switch from the network, you can map it to a forwarding class and a loss priority based on its DSCP code points by applying a classifier to the Layer 3 physical interface. The forwarding class is mapped to an egress queue and to scheduling properties.

By default, before a packet exits the network-facing interface on the ingress switch, the switch copies the DSCP code points from the packet header into the VXLAN header, so the DSCP code points are not rewritten. However, you can configure a rewrite rule on the egress interface (network-facing interface) of the ingress switch if you want to change the value of the DSCP code points.

On the egress switch, the network-facing interface reads the DSCP code points from the VXLAN header and assigns packets to forwarding classes (which are mapped to egress queues) and loss priorities based on the DSCP code points.

NOTE: You cannot classify traffic using an IEEE 802.1p classifier.

Rewrite Rules on OVSDB-Managed VXLAN Interfaces

When packets exit a network, edge switches might need to change the CoS settings of the packets. Rewrite rules change the value of the code points in the packet header by rewriting the code points to a

different value in the outgoing packet. See ["Understanding CoS Rewrite Rules" on page 125](#) for detailed information about rewrite rules.

On OVSDB-managed VXLAN interfaces, you can apply DSCP rewrite rules to packets on network-facing physical interfaces. You cannot apply rewrite rules to access-facing OVSDB-managed VXLAN interfaces, and you cannot apply rewrite rules to IEEE 802.1p code points on network-facing interfaces.

By default, before a packet exits the network-facing interface on the ingress switch, the switch copies the DSCP code points from the packet header into the VXLAN header, so the DSCP code points are not rewritten. The VXLAN header needs to contain the correct DSCP code points because the network-facing ingress port of the egress switch uses the DSCP code points in the VXLAN header to classify the incoming packets.

If you want to change the value of the DSCP code points before the switch transmits packets across the network to the egress switch, you can configure a DSCP rewrite rule and apply it to the egress (network-facing) interface on the ingress switch.

NOTE: Rewrite rules on OVSDB-managed VXLAN interfaces rewrite only the DSCP code point value in the VXLAN header. Rewrite rules on OVSDB-managed VXLAN interfaces do not rewrite the inner (IP) packet header DSCP code point value, so the DSCP code point value in the IP packet header remains unchanged.

Schedulers on OVSDB-Managed VXLAN Interfaces

Packet scheduling (the allocation of port resources such as bandwidth, scheduling priority, and buffers) on OVSDB-managed VXLAN interfaces uses enhanced transmission selection (ETS) hierarchical port scheduling, the same as other interfaces on the switch.

ETS hierarchical port scheduling allocates port bandwidth to traffic in two tiers. ETS provides better port bandwidth utilization and greater flexibility to allocate port resources to forwarding classes (this equates to allocating port resources to output queues because queues are mapped to forwarding classes) and to groups of forwarding classes called forwarding class sets (fc-sets).

First, ETS allocates port bandwidth to fc-sets (also known as priority groups). Each fc-set consists of one or more forwarding classes that carry traffic that requires similar CoS treatment. The bandwidth each fc-set receives is then allocated to the forwarding classes in that fc-set. Each forwarding class is mapped to an output queue. The scheduling properties of a forwarding class are assigned to the queue to which the forwarding class is mapped. Traffic control profiles control the allocation of port bandwidth to fc-sets. Queue schedulers control the allocation of fc-set bandwidth to forwarding classes. See ["Understanding CoS Output Queue Schedulers" on page 338](#), ["Understanding CoS Traffic Control Profiles" on page 401](#), and ["Understanding CoS Hierarchical Port Scheduling \(ETS\)" on page 438](#) for detailed information about scheduling.

NOTE: It is important to take into account the overhead due to VXLAN header encapsulation when you calculate the amount of bandwidth to allocate to VXLAN traffic. When a virtual tunnel endpoint (VTEP) encapsulates a packet in VXLAN, the VXLAN header adds 50 bytes to the packet.

When you configure the queue scheduler transmit rate, which is the minimum amount of guaranteed bandwidth allocated to traffic mapped to a particular queue, and the traffic control profile guaranteed rate, which is the minimum amount of guaranteed bandwidth allocated to traffic mapped to a particular priority group (fc-set), be sure to configure a high enough bandwidth allocation to account for the VXLAN header overhead.

RELATED DOCUMENTATION

[Understanding CoS Classifiers | 96](#)

[Understanding CoS Rewrite Rules | 125](#)

[Understanding CoS Output Queue Schedulers | 338](#)

[Understanding CoS Traffic Control Profiles | 401](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

[Understanding Default CoS Scheduling and Classification | 321](#)

Understanding the OVSDb Protocol Running on Juniper Networks Devices

[Configuring CoS on OVSDb-Managed VXLAN Interfaces | 78](#)

[Example: Configuring Unicast Classifiers | 113](#)

[Defining CoS Rewrite Rules | 128](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

Manually Configuring VXLANs on QFX Series and EX4600 Switches

Configuring CoS on OVSDb-Managed VXLAN Interfaces

On QFX5100 and QFX10000 Series switches, you can configure packet classification, packet scheduling, and packet code point rewrite (rewrite rules) class of service (CoS) features on OVSDb-managed VXLAN interfaces. An OVSDb-managed VXLAN interface uses an OVSDb controller to create and manage the VXLAN interfaces and tunnels.

Classifier, scheduler, and rewrite rule configuration on OVSDB-managed VXLAN interfaces uses the same CLI statements as CoS configuration on regular Ethernet interfaces. However, feature support differs on OVSDB-managed VXLAN interfaces compared to regular Ethernet interfaces in several ways, depending on whether a switch interface is access-facing (connected to servers and other devices accessing the network) or network-facing (connected to the network, for example, switch interfaces that connect to a VXLAN gateway).

- **Classifiers**—On access-facing ingress interfaces, you can configure either BA or MF DSCP classifiers.
On network-facing ingress interfaces, you can configure only DSCP classifiers.
- **Rewrite rules**—On network-facing interfaces, you can configure DSCP rewrite rules. Access-facing interfaces do not support rewrite rules. IEEE 802.1p rewrite rules are not supported.

NOTE: Rewrite rules rewrite the DSCP code point on the VXLAN header only. Rewrite rules do not rewrite the DSCP code point on the inner packet header. If you do not configure a rewrite rule, by default, the code point value in the packet header is copied into the VXLAN header.

- **Schedulers**—Egress interfaces use enhanced transmission selection (ETS) hierarchical port scheduling, the same as regular Ethernet interfaces, and the same features are supported. You can configure packet scheduling on access-facing and network-facing egress interfaces.

For more information about CoS feature support on OVSDB-managed VXLAN interfaces, (see ["Understanding CoS on OVSDB-Managed VXLAN Interfaces" on page 73.](#))

NOTE: This topic covers CoS configuration on OVSDB-managed VXLAN interfaces. It does not cover OVSDB or VXLAN configuration. See *Understanding Dynamically Configured VXLANs in an OVSDB Environment* for information about OVSDB-managed VXLANs.

NOTE: If you do not configure CoS on an interface, the interface uses the default CoS properties. If you configure some CoS properties on an interface, the interface uses the configured CoS for those properties and default CoS for unconfigured properties. The only difference in the default settings on OVSDB-managed VXLAN interfaces is that if you do not configure a rewrite rule, by default, the code point value in the packet header is copied into the VXLAN header. (There is no default rewrite rule on other interfaces.) See ["Understanding Default CoS Scheduling and Classification" on page 321](#) for information about default scheduler and classifier settings.

The following three procedures show how to configure classifiers, rewrite rules, and ETS hierarchical port scheduling on OVSDB-managed VXLAN interfaces.

You can configure classifiers based on the default classifier or a previously configured classifier, or you can create completely new classifiers that do not use any default values. This example is for a network interface.

1. To configure a classifier on an ingress interface using the default classifier or a previously configured classifier as a template (the switch uses the default values for any values that you do not explicitly configure), include the `import` statement and specify `default` or the classifier name as the classifier to import, and associate the classifier with a forwarding class, a loss priority, and one or more code points:

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

To create a classifier that is not based on the default classifier or a previously existing classifier, create a new classifier and associate it with a forwarding class, a loss priority, and one or more code points:

```
[edit class-of-service classifiers]
user@switch# set (dscp | c-dscp) forwarding-class forwarding-class-name loss-priority level low-code-point code
```

NOTE: On network-facing ingress interfaces, only BA DSCP classifiers are supported. Access-facing ingress interfaces support both BA and MF DSCP classification.

2. Apply the classifier to one or more OVSDB-managed VXLAN interfaces on the switch:

```
[edit class-of-service interfaces]
user@switch# set interfaces interface-name classifiers dscp | c-dscp
```

You can configure rewrite rules based on the default rewrite rule or a previously existing rewrite rule. The default rewrite rule writes the inner packet header value to the VXLAN outer header, or you can create completely new classifiers that do not use any default values. You can configure rewrite rules only on network-facing interfaces, and the only supported rewrite rules are DSCP rewrite rules.

1. To configure a rewrite rule on a network-facing egress interface using the default rewrite rule or a previously configured rewrite rule as a template (the switch uses the default values for any values that you do not explicitly configure), include the `import` statement and specify `default` or the rewrite

rule name as the rewrite rule to import, and associate the rewrite rule with a forwarding class, a loss priority, and one or more code points:

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

To create a rewrite rule that is not based on the default rewrite rule or a previously existing rewrite rule, create a new rewrite rule and associate it with a forwarding class, a loss priority, and one or more code points:

```
[edit class-of-service rewrite-rules]
user@switch# set dscp rewrite-name forwarding-class forwarding-class-name loss-priority level
code-points [aliases] [bit-patterns]
```

NOTE: Rewrite rules are not supported on access-facing interfaces.

2. Apply the rewrite rule to one or more OVSDB-managed VXLAN interfaces on the switch:

```
[edit class-of-service interfaces]
user@switch# set interface-name unit unit rewrite-rules dscp rewrite-name
```

ETS hierarchical port scheduling allocates port bandwidth to traffic in two tiers. ETS provides better port bandwidth utilization and greater flexibility to allocate port resources to forwarding classes (this equates to allocating port resources to output queues because queues are mapped to forwarding classes) and to groups of forwarding classes called forwarding class sets (fc-sets).

First, ETS allocates port bandwidth to fc-sets (also known as priority groups). Each fc-set consists of one or more forwarding classes that carry traffic that requires similar CoS treatment. The bandwidth each fc-set receives is then allocated to the forwarding classes in that fc-set. Each forwarding class is mapped to an output queue. The scheduling properties of a forwarding class are assigned to the queue to which the forwarding class is mapped. Traffic control profiles control the allocation of port bandwidth to fc-sets. Queue schedulers control the allocation of fc-set bandwidth to forwarding classes. See ["Understanding CoS Output Queue Schedulers" on page 338](#), ["Understanding CoS Traffic Control Profiles" on page 401](#), and ["Understanding CoS Hierarchical Port Scheduling \(ETS\)" on page 438](#) for detailed information about scheduling.

Schedulers define the CoS properties of the output queues mapped to forwarding classes. After you configure a scheduler, you use a scheduler map to map the scheduler to one or more forwarding classes. Mapping the scheduler to a forwarding class applies the scheduling properties to the traffic in the forwarding class.

Schedulers define the following characteristics for the forwarding classes (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 an fc-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 weighted random early detection (WRED) packet drop characteristics to the scheduler.

NOTE: If ingress port congestion occurs because of egress port congestion, apply a drop profile to the traffic on the congested egress port so that traffic is dropped at the egress interface instead of at the ingress interface. (Ingress interface congestion can affect uncongested ports when an ingress port transmits traffic to both congested and uncongested egress ports.)

- **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.

A traffic control profile defines the CoS properties of an fc-set, and the amount of port resources allocated to the group of forwarding classes (queues) in the fc-set. After you configure a traffic control profile, you apply it (with an associated fc-set) to an interface, to configure scheduling on that interface for traffic that belongs to the forwarding classes in the fc-set .

A traffic control profile defines the following characteristics for the fc-set (priority group) mapped to the traffic control profile when you apply traffic control profile and fc-set to an interface:

- **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 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 in each fc-set.

The following procedure shows how to configure scheduler properties, map schedulers to forwarding classes, map forwarding classes to fc-sets, configure traffic control profile properties, and apply traffic control profiles and fc-sets to interfaces (to apply the ETS ports scheduling configuration to interfaces).

NOTE: You do not have to explicitly configure all of the scheduler and traffic control profile characteristics. Some characteristics are disabled by default, such as ECN, and should only be enabled under certain conditions. You can have a mix of configured CoS properties and default CoS properties.

1. Name the queue scheduler and define the minimum guaranteed bandwidth for the queue:

```
[edit class-of-service]
user@switch# set schedulers scheduler-name transmit-rate (rate | percent
percentage)
```

2. Define the maximum bandwidth for the queue:

```
[edit class-of-service schedulers scheduler-name]
user@switch# set shaping-rate (rate | percent percentage)
```

3. Define the queue priority:

```
[edit class-of-service schedulers scheduler-name]
user@switch# set priority level
```

4. Define the drop profile 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 20
```

6. Enable ECN, if desired (queue should handle best-effort traffic):

```
[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
```

This completes the characteristics you can configure in a scheduler, and scheduler mapping to forwarding classes. The next steps show how to configure traffic control profiles.

8. 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)
```

9. 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)
```

10. Attach a scheduler map to the traffic control profile; the scheduler map associates the schedulers and forwarding classes (queues) in the scheduler map with the traffic control profile:

```
[edit class-of-service traffic-control-profiles traffic-control-profile-name]
user@switch# set scheduler-map scheduler-map-name
```

This completes the characteristics you can configure in a traffic control profile. The next step shows how to assign forwarding classes to fc-sets.

11. Assign one or more forwarding classes to the fc-set:

```
[edit class-of-service]
user@switch# set forwarding-class-sets forwarding-class-set-name class forwarding-class-name
```

This completes assigning forwarding classes to fc-sets. The next steps show how to apply ETS hierarchical port scheduling to interfaces.

12. To apply ETS hierarchical port scheduling to interfaces, associate an fc-set and a traffic control profile with interfaces. The fc-set determines the forwarding class(es) and queue(s) that use the specified interface. The traffic control profile determines the amount of port resources allocated to the fc-set, and the mapping of forwarding classes to schedulers in the traffic control profile determines the allocation of fc-set resources to the forwarding classes that are members of the fc-set.

```
user@switch# set interfaces interface-name forwarding-class-set fc-set-name output-traffic-control-profile tcp-name
```

RELATED DOCUMENTATION

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[Example: Configuring Unicast Classifiers | 113](#)

[Configuring CoS WRED Drop Profiles | 284](#)

[Example: Configuring ECN | 307](#)

[Understanding CoS on OVSDb-Managed VXLAN Interfaces | 73](#)

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](#) | 21

CoS Code-Point Aliases

IN THIS CHAPTER

- Understanding CoS Code-Point Aliases | 90
- Defining CoS Code-Point Aliases | 93
- Monitoring CoS Code-Point Value Aliases | 94

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](#)" on page 90.

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 28 on page 91](#) shows the default mapping of code-point aliases to IEEE code points.

Table 28: 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 29 on page 91](#) shows the default mapping of code-point aliases to DSCP and DSCP IPv6 code points.

Table 29: Default DSCP and DSCP IPv6 Code-Point Aliases

CoS Value Types	Mapping
ef	101110
af11	001010
af12	001100
af13	001110

Table 29: Default DSCP and DSCP IPv6 Code-Point Aliases *(Continued)*

CoS Value Types	Mapping
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
cs4	100000
cs5	101000

Table 29: Default DSCP and DSCP IPv6 Code-Point Aliases (*Continued*)

CoS Value Types	Mapping
nc1	110000
nc2	111000

RELATED DOCUMENTATION

[Understanding Junos CoS Components | 21](#)

[Defining CoS Code-Point Aliases | 93](#)

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 | 94](#)

[Understanding CoS Code-Point Aliases | 90](#)

Monitoring CoS Code-Point Value Aliases

IN THIS SECTION

- [Purpose | 94](#)
- [Action | 94](#)
- [Meaning | 95](#)

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 30 on page 95 summarizes key output fields for CoS value aliases.

Table 30: Summary of Key CoS Value Alias Output Fields

Field	Values
Code point type	Type of the CoS value: <ul style="list-style-type: none">dscp—Examines Layer 3 packet headers for IP packet classification.dscp-ipv6—Examines Layer 3 packet headers for IPv6 packet classification.ieee-802.1—Examines Layer 2 packet headers for packet classification.exp—Examines MPLS packet headers for packet classification. <p>NOTE: 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

CHAPTER 4

CoS Classifiers

IN THIS CHAPTER

- Understanding CoS Classifiers | 96
- Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p) | 106
- Example: Configuring Classifiers | 108
- Example: Configuring Unicast Classifiers | 113
- Example: Configuring Multidestination (Multicast, Broadcast, DLF) Classifiers | 117
- Understanding Host Inbound Traffic Classification | 121
- Configuring a Global MPLS EXP Classifier | 122
- Monitoring CoS Classifiers | 123

Understanding CoS Classifiers

IN THIS SECTION

- Interfaces and Output Queues | 97
- Output Queues for Unicast and Multidestination Traffic | 98
- Classifier Support by Type | 99
- Behavior Aggregate Classifiers | 100
- Fixed Classifiers on Ethernet Interfaces | 103
- Fixed Classifiers on Native Fibre Channel Interfaces (NP_Ports) | 104
- Multifield Classifiers | 105
- MPLS EXP Classifiers | 105
- Packet Classification for IRB Interfaces and RVIs | 105

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 multdestination (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 multdestination traffic. You cannot mix unicast and multdestination traffic in the same classifier or forwarding class. On QFX10000 switches, a classifier can assign both unicast and multdestination 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](#)" on page 130 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 multdestination traffic queues. You can apply only multdestination BA classifiers to multdestination queues. A multdestination BA classifier should contain only forwarding classes that are mapped to multdestination queues.

You can apply unicast classifiers to one or more interfaces. Multdestination classifiers and EXP classifiers apply to all of the switch interfaces and cannot be applied to individual interfaces. Use the DSCP multdestination classifier for both IP and IPv6 multdestination traffic. The DSCP IPv6 classifier is not supported for multdestination 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 31 on page 99](#) shows how many of each type of classifiers you can configure, and how many entries you can configure per classifier.

Table 31: 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 32 on page 100](#) lists different types of interfaces and the corresponding implicit default BA classifiers.

Table 32: 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

Table 32: Default BA Classification (*Continued*)

Type of Interface	Default BA Classification
(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 multidestination classifier is the IEEE 802.1p multidestination 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 33 on page 103](#). 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 33: 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 | 26](#)

[Understanding Default CoS Settings | 30](#)

[Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 106](#)

[Example: Configuring Unicast Classifiers | 113](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 106](#)

[Example: Configuring Multidestination \(Multicast, Broadcast, DLF\) Classifiers | 117](#)

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 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 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 multdestination traffic, except on QFX10000 switches or NFX Series devices, configure the classifier as a multdestination 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 \(ETS\) | 445](#)

[Example: Configuring Unicast Classifiers | 113](#)

Configuring a Global MPLS EXP Classifier

Configuring Rewrite Rules for MPLS EXP Classifiers

[Monitoring CoS Classifiers | 123](#)

[Understanding CoS Classifiers | 96](#)

[Understanding CoS Classifiers](#)

Understanding CoS MPLS EXP Classifiers and Rewrite Rules

[Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130](#)

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

Example: Configuring Classifiers

IN THIS SECTION

- [Requirements | 110](#)
- [Overview | 110](#)
- [Verification | 111](#)

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 Classifiers

Step-by-Step Procedure

To configure an IEEE 802.1 BA classifier named `ba-classifier` as the default IEEE 802.1 classifier:

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-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-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-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-classifier forwarding-class nc loss-priority low code-points
110
```

5. Apply the 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-classifier
```

Requirements

This example uses the following hardware and software components:

- One switch.
- Junos OS Release 15.1X53-D10 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 Differentiated Services code point (DSCP or DSCP IPv6) value, IEEE 802.1p value, or MPLS EXP value. (EXP classifiers can be applied only to 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 VLAN header or the DSCP bits in the IP 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.

This example describes how to configure a BA classifier called `ba-classifier` as the default IEEE 802.1 mapping of incoming traffic to forwarding classes, and apply it to ingress interface `xe-0/0/10`. The BA classifier assigns loss priorities, as shown in [Table 34 on page 111](#), to incoming packets in the four default forwarding classes. You can adapt the example to DSCP traffic by specifying a DSCP classifier instead of an IEEE classifier, and by applying DSCP bits instead of CoS bits.

To set multifield classifiers, use firewall filter rules.

Table 34: ba-classifier Loss Priority Assignments

Forwarding Class	CoS Traffic Type	ba-classifier Loss Priority to IEEE 802.1p Code Point Mapping	Packet Drop Attribute
be	Best-effort traffic	Low loss priority code point: 000	drop
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	no-loss
nc	Network-control traffic	Low loss priority code point: 110	drop

Verification

IN THIS SECTION

- [Verifying the Classifier Configuration | 111](#)
- [Verifying the Ingress Interface Configuration | 112](#)

To verify the classifier configuration, perform these tasks:

Verifying the Classifier Configuration

Purpose

Verify that you configured the 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-classifier`:

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

RELATED DOCUMENTATION

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 106](#)

[Configuring a Global MPLS EXP Classifier](#)

[Configuring Rewrite Rules for MPLS EXP Classifiers](#)

[Monitoring CoS Classifiers | 123](#)

[Understanding CoS Classifiers | 96](#)

[Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130](#)

Example: Configuring Unicast Classifiers

IN THIS SECTION

- [Requirements | 114](#)
- [Overview | 114](#)
- [Verification | 115](#)

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 multdestination (multicast, broadcast, and destination lookup fail) traffic to different classifiers. One classifier cannot include both unicast and multdestination 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 35 on page 115](#), 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 35: 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
nc	Network-control traffic	Low loss priority code point: 110	drop

Verification

IN THIS SECTION

● [Verifying the Unicast Classifier Configuration](#) | 116

- [Verifying the Ingress Interface Configuration | 116](#)

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 \(ETS\) | 445](#)

[Example: Configuring Multidestination \(Multicast, Broadcast, DLF\) Classifiers | 117](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 106](#)

[*Configuring a Global MPLS EXP Classifier*](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 106](#)

[Monitoring CoS Classifiers | 123](#)

[Understanding CoS Classifiers | 96](#)

[Understanding CoS Classifiers](#)

[Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130](#)

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

IN THIS SECTION

- [Requirements | 118](#)
- [Overview | 119](#)
- [Verification | 120](#)

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 multdestination classifiers. On EX4300 switches, you can apply multdestination classifiers globally or to a specific interface. If you apply multdestination classifiers both globally and to a specific interface, the classifications on the interface take precedence.

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, except on an EX4300 switch.

Unicast and multdestination traffic must use different classifiers.

Configuring Multdestination 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 36 on page 119](#), to incoming packets in the multidestination forwarding class.

You can also use firewall filters to set multifield classifiers.

Table 36: 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 | 120](#)
- [Verifying the Multidestination Classifier Configuration | 120](#)

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

Example: Configuring Unicast Classifiers 113
Defining CoS BA Classifiers (DSCP, DSCP IPv6, IEEE 802.1p) 106
Monitoring CoS Classifiers 123
Understanding CoS Classifiers 96
Understanding CoS Classifiers
Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces 130

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.

Configuring a Global MPLS EXP Classifier

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

You can configure up to 64 EXP classifiers, however, the switch uses only one MPLS EXP classifier as a global classifier, which is applied only on interfaces configured as `family mpls`. All `family mpls` switch interfaces use the global EXP classifier to classify MPLS traffic.

There is no default EXP classifier. If you want to classify incoming MPLS packets using the EXP bits, you must configure a global EXP classifier. The global classifier applies to all MPLS traffic on all `family mpls` interfaces.

If a global EXP classifier is configured, MPLS traffic on `family mpls` interfaces uses the EXP classifier. 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. DSCP classifiers are not applied to MPLS traffic.

To configure an MPLS EXP classifier using the CLI:

1. Create an EXP classifier and associate it with a forwarding class, a loss priority, and a code point:

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

2. Apply the EXP classifier to the switch interfaces:

```
[edit class-of-service]
user@switch# set system-defaults classifiers exp classifier-name
```

Monitoring CoS Classifiers

IN THIS SECTION

- Purpose | 123
- Action | 123
- Meaning | 123

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 37 on page 123](#) summarizes key output fields for CoS classifiers.

Table 37: Summary of Key CoS Classifier Output Fields

Field	Values
Classifier	Name of a classifier.

Table 37: Summary of Key CoS Classifier Output Fields (*Continued*)

Field	Values
Code point type	<p>Type of classifier:</p> <ul style="list-style-type: none"> dscp—All classifiers of the DSCP type. ieee-802.1—All classifiers of the IEEE 802.1 type. ieee-mcast—All classifiers of the IEEE 802.1 multicast type. <p>NOTE: 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"> exp—All classifiers of the MPLS exp type. <p>NOTE: 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.

CoS Rewrite Rules

IN THIS CHAPTER

- Understanding CoS Rewrite Rules | 125
- Defining CoS Rewrite Rules | 128
- Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130
- Troubleshooting an Unexpected Rewrite Value | 145
- Understanding CoS MPLS EXP Classifiers and Rewrite Rules | 147
- Configuring Rewrite Rules for MPLS EXP Classifiers | 151
- Monitoring CoS Rewrite Rules | 153

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 38 on page 127](#) shows how many of each type of rewrite rules you can configure, and how many entries you can configure per rewrite rule.

Table 38: 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
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 | 21](#)

[Defining CoS Rewrite Rules | 128](#)

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 153
Configuring Rewrite Rules for MPLS EXP Classifiers
Understanding CoS Rewrite Rules 125
Understanding CoS MPLS EXP Classifiers and Rewrite Rules

Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces

IN THIS SECTION

- Supported Classifier and Rewrite Rule Types | 130
- Ethernet Interfaces Supported for Classifier and Rewrite Rule Configuration | 133
- Default Classifiers | 137
- Default Rewrite Rules | 137
- Classifier Precedence | 137
- Classifier Behavior and Limitations | 139
- Rewrite Rule Precedence and Behavior | 140
- Classifier and Rewrite Rule Configuration Interaction with Ethernet Interface Configuration | 141

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 39 on page 131](#) shows the supported types of classifiers and rewrite rules supports:

Table 39: 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.
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>

Table 39: Supported Classifiers and Rewrite Rules (Continued)

Classifier or Rewrite Rule Type	Description
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.
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*).

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*.

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 40 on page 134](#) 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 40: 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

Table 40: 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
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.			
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 multidestination and DSCP multidestination classifiers are applied to all interfaces and cannot be applied to individual interfaces. No DSCP IPv6 multidestination classifier is supported. IPv6 multidestination traffic uses the DSCP multidestination 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 41 on page 136](#) 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 41: 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 multdestination traffic, the default multdestination classifier is the IEEE 802.1p multdestination 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 \(DSCP, DSCP IPv6, IEEE 802.1p\)" on page 106](#) 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 | 26](#)

[Configuring CoS | 14](#)

Troubleshooting an Unexpected Rewrite Value

IN THIS SECTION

- Problem | 145
- Cause | 145
- Solution | 146

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](#)

Understanding CoS MPLS EXP Classifiers and Rewrite Rules

IN THIS SECTION

- [EXP Classifiers | 148](#)
- [EXP Rewrite Rules | 150](#)
- [Schedulers | 151](#)

You can use *class of service* (CoS) within MPLS networks to prioritize certain types of traffic during periods of congestion by applying packet classifiers and rewrite rules to the MPLS traffic. MPLS classifiers are global and apply to all interfaces configured as `family mpls` interfaces.

When a packet enters a customer-edge interface on the ingress provider edge (PE) switch, the switch associates the packet with a particular CoS servicing level before placing the packet onto the label-switched path (LSP). The switches within the LSP utilize the CoS value set at the ingress PE switch to determine the CoS service level. The CoS value embedded in the classifier is translated and encoded in the MPLS header by means of the experimental (EXP) bits.

EXP classifiers map incoming MPLS packets to a forwarding class and a loss priority, and assign MPLS packets to output queues based on the forwarding class mapping. EXP classifiers are behavior aggregate (BA) classifiers.

EXP rewrite rules change (rewrite) the CoS value of the EXP bits in outgoing packets on the egress queues of the switch so that the new (rewritten) value 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.

NOTE: On QFX5200, QFX5100, QFX3500, QFX3600, and EX4600 switches, and on QFabric systems, there is no default EXP classifier. If you want to classify incoming MPLS packets using the EXP bits, you must configure a global EXP classifier. The global EXP classifier applies to all MPLS traffic on interfaces configured as `family mpls`.

On QFX10000 switches, there is a no default EXP classifier. If you want to classify incoming MPLS packets using the EXP bits, you must configure EXP classifiers and apply them to logical interfaces configured as `family mpls`. (You cannot apply classifiers to physical interfaces.). You can configure up to 64 EXP classifiers.

There is no default EXP rewrite rule. If you want to rewrite the EXP bit value at the egress interface, you must configure EXP rewrite rules and apply them to logical interfaces.

EXP classifiers and rewrite rules are applied only to interfaces that are configured as `family mpls` (for example, set interfaces `xe-0/0/35 unit 0 family mpls`.)

This topic includes:

EXP Classifiers

On QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, unlike DSCP and IEEE 802.1p BA classifiers, EXP classifiers are global to the switch and apply to all switch 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.

When you configure and apply an EXP classifier, MPLS traffic on all `family mpls` interfaces uses the EXP classifier, even on interfaces that also have a fixed classifier. If an interface has both an EXP classifier and a fixed classifier, the EXP classifier is applied to MPLS traffic and the fixed classifier is applied to all other traffic.

Also unlike DSCP and IEEE 802.1p BA classifiers, there is no default EXP classifier. If you want to classify MPLS traffic based on the EXP bits, you must explicitly configure an EXP classifier and apply it to the switch interfaces. Each EXP classifier has eight entries that correspond to the eight EXP CoS values (0 through 7, which correspond to CoS bits 000 through 111).

You can configure up to 64 EXP classifiers.

However, on QFX5200, QFX5100, EX4600, and legacy CLI switches, 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 that classifier as the global EXP classifier by including the EXP classifier in the [edit class-of-service system-defaults classifiers exp] hierarchy level. All switch interfaces configured as family mpls use the global EXP classifier to classify MPLS traffic.

On these switches, only one EXP classifier can be configured as the global EXP classifier at any time. If you want to change the global EXP classifier, delete the global EXP classifier configuration (use the **user@switch# delete class-of-service system-defaults classifiers exp** configuration statement), then configure the new global EXP classifier.

NOTE: QFX5130 switch does not support MPLS CoS.

QFX10000 switches do not support global EXP classifiers. You can configure one EXP classifier and apply it to multiple logical interfaces, or configure multiple EXP classifiers and apply different EXP classifiers to different logical interfaces.

If an EXP classifier is not configured, then if a fixed classifier is applied to the interface, the MPLS traffic uses the fixed classifier. (Switches that have a default EXP classifier use the default classifier.) If no EXP classifier and no fixed classifier are applied to the interface, MPLS traffic is treated as best-effort traffic using the 802.1 default untrusted classifier. DSCP classifiers are not applied to MPLS traffic.

On QFX5200, QFX5100, EX4600, and legacy CLI switches, 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.

NOTE: The switch uses only the outermost label of incoming EXP packets for classification.

NOTE: MPLS packets with 802.1Q tags are not supported.

On QFX5220 switch, you can use class of service (CoS) within MPLS networks to prioritize certain types of traffic during periods of congestion by applying packet classifiers and rewrite rules to the MPLS traffic. We have also added the MPLS EXP rewrite support.

- Default CoS on the Provider (P) and Provider Edge (PE) routers for MPLS interfaces – The MPLS traffic uses the default EXP classifier. MPLS traffic is treated as best-effort traffic using the 802.1 default untrusted classifier. The default EXP classifier applies to all MPLS traffic on interfaces configured as family mpls. DSCP classifiers are not applied to MPLS traffic.
- Default CoS on PE routers for Layer 3 interfaces – By default, all L3VPN logical interfaces are bound to default Differentiated Services Code Point (DSCP) classifiers.

If you apply an EXP classifier on a penultimate hop popping (PHP) node, then by default, the IP header time-to-live (TTL) value is overwritten by the MPLS header TLL value, and the IP header DSCP bits are over written by a zero (0), which signifies uniform mode. On Junos OS Evolved, to use pipe mode, where IP header TTL and IP header DSCP bits are not overwritten, you should configure the following command:

```
set protocols mpls no-propagate-ttl
```

However, on Junos OS, you can configure MPLS CoS without the `set protocols mpls no-propagate-ttl` command.

NOTE: The DSCP of IP in MPLS packets can't be remarked either at PE or P routers.

EXP Rewrite Rules

As MPLS packets enter or exit a network, edge switches might be required to alter the class-of-service (CoS) settings of the packets. EXP *rewrite rules* set the value of the EXP CoS bits within the header of the outgoing MPLS packet on family `mpls` interfaces. Each rewrite rule reads the current forwarding class and loss priority associated with the packet, locates the chosen CoS value from a table, and writes that CoS value into the packet header, replacing the old CoS value. EXP rewrite rules apply only to MPLS traffic.

EXP rewrite rules apply only to logical interfaces. You cannot apply EXP rewrite rules to physical interfaces.

There are no default EXP rewrite rules. If you want to rewrite the EXP value in MPLS packets, you must configure EXP rewrite rules and apply them to logical interfaces. If no rewrite rules are applied, all MPLS labels that are pushed have a value of zero (0). The EXP value remains unchanged on MPLS labels that are swapped.

You can configure up to 64 EXP rewrite rules, but you can only apply 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.

You can apply an EXP rewrite rule to an interface that has a DSCP, DSCP IPv6, or IEEE 802.1p rewrite rule. Only MPLS traffic uses the EXP rewrite rule. MPLS traffic does not use DSCP or DSCP IPv6 rewrite rules.

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.

Schedulers

The schedulers for using CoS with MPLS are the same as for the other CoS configurations on the switch. Default schedulers are provided only for the best-effort, fcoe, no-loss, and network-control default forwarding classes. If you configure a custom forwarding class for MPLS traffic, you need to configure a scheduler to support that forwarding class and provide bandwidth to that forwarding class.

Configuring Rewrite Rules for MPLS EXP Classifiers

You configure EXP rewrite rules to alter CoS values in outgoing MPLS packets on the outbound family `mpls` interfaces of a switch to match 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 an EXP 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 a logical family `mpls` interface. EXP rewrite rules can only be enabled on logical family `mpls` interfaces, not on physical interfaces or on interfaces of other family types. You can also apply an existing EXP rewrite rule on a logical interface.

NOTE: There are no default rewrite rules.

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 family `mpls` 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: 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.

To create an EXP rewrite rule for MPLS traffic and enable it on a logical interface:

1. Create an EXP rewrite rule:

```
user@switch# set class-of-service rewrite-rules exp rewrite-rule-name forwarding-class
forwarding-class-name loss-priority level code-points [aliases] [bit-patterns]
```

For example, to configure an EXP rewrite rule named `exp-rr-1` for a forwarding class named `mpls-1` with a loss priority of `low` that rewrites the EXP code point value to `001`:

```
user@switch# set class-of-service rewrite-rules exp exp-rr-1 forwarding-class mpls-1 loss-
priority low code-points 001
```

2. Apply the rewrite rule to a logical interface:

```
user@switch # set class-of-service interfaces interface-name unit logical-unit rewrite-rules
exp rewrite-rule-name
```

For example, to apply a rewrite rule named `exp-rr-1` to logical interface `xe-0/0/10.0`:

```
user@switch# set class-of-service interfaces xe-0/0/10 unit 0 rewrite-rules exp exp-rr-1
```

NOTE: In this example, all forwarding classes assigned to port `xe-0/0/10` 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 interface.

Monitoring CoS Rewrite Rules

IN THIS SECTION

- Purpose | 153
- Action | 153
- Meaning | 153

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 42 on page 153](#) summarizes key output fields for CoS rewrite rules.

Table 42: Summary of Key CoS Rewrite Rule Output Fields

Field	Values
Rewrite rule	Name of the rewrite rule.

Table 42: Summary of Key CoS Rewrite Rule Output Fields *(Continued)*

Field	Values
Code point type	Rewrite rule type: <ul style="list-style-type: none"> • dscp—For IPv4 DiffServ traffic. • dscp-ipv6—For IPv6 Diffserv traffic. • ieee-802.1—For Layer 2 traffic. • exp—For MPLS traffic. <p>NOTE: OCX Series switches do not support MPLS.</p>
Index	Internal index for the rewrite rule.
Forwarding class	Name of the forwarding class that is used to determine CoS values for rewriting in combination with loss priority. Rewrite rules are applied to CoS values in outgoing packets based on forwarding class and loss priority setting.
Loss priority	Level of loss priority that is used to determine CoS values for rewriting in combination with forwarding class.
Code point	Rewrite code point value.

RELATED DOCUMENTATION

[Defining CoS Rewrite Rules](#) | 128

CoS Forwarding Classes and Forwarding Class Sets

IN THIS CHAPTER

- [Understanding CoS Forwarding Classes | 155](#)
- [Defining CoS Forwarding Classes | 162](#)
- [Forwarding Policy Options Overview | 164](#)
- [Configuring CoS-Based Forwarding | 166](#)
- [Example: Configuring CoS-Based Forwarding | 170](#)
- [Example: Configuring Forwarding Classes | 174](#)
- [Understanding CoS Forwarding Class Sets \(Priority Groups\) | 181](#)
- [Defining CoS Forwarding Class Sets | 183](#)
- [Example: Configuring Forwarding Class Sets | 184](#)
- [Monitoring CoS Forwarding Classes | 189](#)

Understanding CoS Forwarding Classes

IN THIS SECTION

- [Default Forwarding Classes | 157](#)
- [Forwarding Class Configuration Rules | 158](#)
- [Lossless Transport Support | 160](#)

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 43 on page 157 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 43: 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 43: 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 44 on page 158](#) applies only to multidestination traffic except on QFX10000 switches and NFX Series devices.

Table 44: 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* 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

[Understanding Junos CoS Components | 21](#)

[Understanding CoS Packet Flow | 26](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

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 \(ETS\) | 445](#)

[Example: Configuring Forwarding Classes | 174](#)

[Monitoring CoS Forwarding Classes | 189](#)

[Understanding CoS Forwarding Classes | 155](#)

[Understanding CoS Port Schedulers | 368](#)

Forwarding Policy Options Overview

Class-of-service (CoS)-based forwarding (CBF) enables you to control next-hop selection based on a packet's *class of service* and, in particular, the value of the IP packet's precedence bits.

For example, you might want to specify a particular interface or next hop to carry high-priority traffic while all best-effort traffic takes some other path. When a routing protocol discovers equal-cost paths,

Junos picks a path at random or load-balance across the paths through either hash selection or round robin. CBF allows path selection based on class.

To configure CBF properties, include the following statements at the [edit class-of-service] hierarchy level:

```
[edit class-of-service]
forwarding-policy {
  next-hop-map map-name {
    forwarding-class class-name {
      next-hop [ next-hop-name ];
      lsp-next-hop [ lsp-regular-expression ];
      non-lsp-next-hop;
      discard;
    }
    forwarding-class-default {
      discard;
      lsp-next-hop [ lsp-regular-expression ];
      next-hop [ next-hop-name ];
      non-lsp-next-hop;
    }
  }
}
class class-name {
  classification-override {
    forwarding-class class-name;
  }
}
}
```

NOTE: Beginning with Junos OS Release 17.1R1, QFX10000 Series switches support CoS-based forwarding. [set class-of-service forwarding-policy class] is not supported on QFX10000 Series switches.

Beginning with Junos OS Release 17.2, MX routers with MPCs or MS-DPCs, VMX, PTX3000 routers, PTX5000 routers, and VPTX support configuring CoS-based forwarding (CBF) for up to 16 forwarding classes. All other platforms support CBF for up to 8 forwarding classes. To support up to 16 forwarding classes for CBF on MX routers, enable enhanced-ip at the [edit chassis network-services] hierarchy level. Enabling enhanced-ip is not necessary on PTX routers to support 16 forwarding classes for CBF.

Release History Table

Release	Description
17.2R1	Beginning with Junos OS Release 17.2, MX routers with MPCs or MS-DPCs, VMX, PTX3000 routers, PTX5000 routers, and VPTX support configuring CoS-based forwarding (CBF) for up to 16 forwarding classes.
17.1R1	Beginning with Junos OS Release 17.1R1, QFX10000 Series switches support CoS-based forwarding. [set class-of-service forwarding-policy class] is not supported on QFX10000 Series switches.

RELATED DOCUMENTATION

Configuring CoS-Based Forwarding

Example: Configuring CoS-Based Forwarding

Configuring CoS-Based Forwarding

You can apply CoS-based forwarding (CBF) only to a defined set of routes. Therefore, you must configure a policy statement as in the following example:

```
[edit policy-options]
policy-statement my-cos-forwarding {
  from {
    route-filter destination-prefix match-type;
  }
  then {
    cos-next-hop-map map-name;
  }
}
```

This configuration specifies that routes matching the route filter are subject to the CoS next-hop mapping specified by *map-name*. For more information about configuring policy statements, see the [Routing Policies, Firewall Filters, and Traffic Policers User Guide](#).

NOTE: On M Series routers (except the M120 and M320 routers), forwarding-class-based matching and CBF do not work as expected if the forwarding class has been set with a multifield filter on an input interface.

Beginning with Junos OS Release 17.2, MX routers with MPCs or MS-DPCs, VMX, PTX3000 routers, and PTX5000 routers support configuring CoS-based forwarding (CBF) for up to 16 forwarding classes. All other platforms support CBF for up to 8 forwarding classes. To support up to 16 forwarding classes for CBF on MX routers, enable `enhanced-ip` at the `[edit chassis network-services]` hierarchy level.

You can configure CBF on a device with the supported number or fewer forwarding classes plus a default forwarding class only. Under this condition, the forwarding class to queue mapping can be either one-to-one or one-to-many. However, you cannot configure CBF when the number of forwarding classes configured exceeds the supported number. Similarly, with CBF configured, you cannot configure more than the supported number of forwarding classes plus a default forwarding class.

To specify a CoS next-hop map, include the `forwarding-policy` statement at the `[edit class-of-service]` hierarchy level:

```
[edit class-of-service]
forwarding-policy {
  next-hop-map map-name {
    forwarding-class class-name {
      discard;
      lsp-next-hop [ lsp-regular-expression ];
      next-hop [ next-hop-name ];
      non-lsp-next-hop;
    }
    forwarding-class-default {
      discard;
      lsp-next-hop [ lsp-regular-expression ];
      next-hop [ next-hop-name ];
      non-lsp-next-hop;
    }
  }
}
```

When you configure CBF with OSPF as the interior gateway protocol (IGP), you must specify the next hop as an interface name or next-hop alias, not as an IPv4 or IPv6 address. This is true because OSPF

adds routes with the interface as the next hop for point-to-point interfaces; the next hop does not contain the IP address. For an example configuration, see *Example: Configuring CoS-Based Forwarding*.

For Layer 3 VPNs, when you use class-based forwarding for the routes received from the far-end provider edge (PE) router within a VRF instance, the software can match the routes based on the attributes that come with the received route only. In other words, the matching can be based on the route within RIB-in. In this case, the route-filter statement you include at the [edit policy-options policy-statement my-cos-forwarding from] hierarchy level has no effect because the policy checks the `bgp.l3vpn.0` table, not the `vrf.inet.0` table.

Junos OS applies the CoS next-hop map to the set of next hops previously defined; the next hops themselves can be located across any outgoing interfaces on the routing device. For example, the following configuration associates a set of forwarding classes and next-hop identifiers:

```
[edit class-of-service forwarding-policy]
next-hop-map map1 {
  forwarding-class expedited-forwarding {
    next-hop next-hop1;
    next-hop next-hop2;
  }
  forwarding-class best-effort {
    next-hop next-hop3;
    lsp-next-hop lsp-next-hop4;
  }
  forwarding-class-default {
    lsp-next-hop lsp-next-hop5;
  }
}
```

In this example, `next-hop N` is either an IP address or an egress interface for some next hop, and `lsp-next-hop N` is a regular expression corresponding to any next hop with that label. Q1 through QN are a set of forwarding classes that map to the specific next hop. That is, when a packet is switched with Q1 through QN, it is forwarded out the interface associated with the associated next hop.

This configuration has the following implications:

- A single forwarding class can map to multiple standard next hops or LSP next hops. This implies that load sharing is done across standard next hops or LSP next hops servicing the same class value. To make this work properly, Junos OS creates a list of the equal-cost next hops and forwards packets according to standard load-sharing rules for that forwarding class.
- If a forwarding class configuration includes LSP next hops and standard next hops, the LSP next hops are preferred over the standard next hops. In the preceding example, if both `next-hop3` and `lsp-next-`

hop4 are valid next hops for a route to which map1 is applied, the forwarding table includes entry lsp-next-hop4 only.

- If next-hop-map does not specify all possible forwarding classes, the default forwarding class is selected as the default. *default-forwarding class* defines the next hop for traffic that does not meet any forwarding class in the next hop map. If the default forwarding class is not specified in the next-hop map, a default is designated randomly. The default forwarding class is the class associated with queue 0.
- For LSP next hops, Junos OS uses UNIX `regex(3)`-style regular expressions. For example, if the following labels exist: `lsp`, `lsp1`, `lsp2`, `lsp3`, the statement `lsp-next-hop lsp` matches `lsp`, `lsp1`, `lsp2`, and `lsp3`. If you do not want this behavior, you must use the anchor characters `lsp-next-hop " ^lsp$"`, which match `lsp` only.
- The route filter does not work because the policy checks against the `bgp.l3vpn.0` table instead of the `vrf.inet.0` table.

The final step is to apply the route filter to routes exported to the forwarding engine. This is shown in the following example:

```

routing-options {
  forwarding-table {
    export my-cos-forwarding;
  }
}

```

This configuration instructs the routing process to insert routes to the forwarding engine matching `my-cos-forwarding` with the associated next-hop CBF rules.

The following algorithm is used when you apply a configuration to a route:

- If the route is a single next-hop route, all traffic goes to that route; that is, no CBF takes effect.
- For each next hop, associate the proper forwarding class. If a next hop appears in the route but not in the `cos-next-hop map`, it does not appear in the forwarding table entry.
- The default forwarding class is used if not all forwarding classes are specified in the next-hop map. If the default is not specified, the default is assigned to the lowest class defined in the next-hop map.

Release History Table

Release	Description
17.2R1	Beginning with Junos OS Release 17.2, MX routers with MPCs or MS-DPCs, VMX, PTX3000 routers, and PTX5000 routers support configuring CoS-based forwarding (CBF) for up to 16 forwarding classes.

RELATED DOCUMENTATION

Load Balancing VPLS Non-Unicast Traffic Across Member Links of an Aggregate Interface

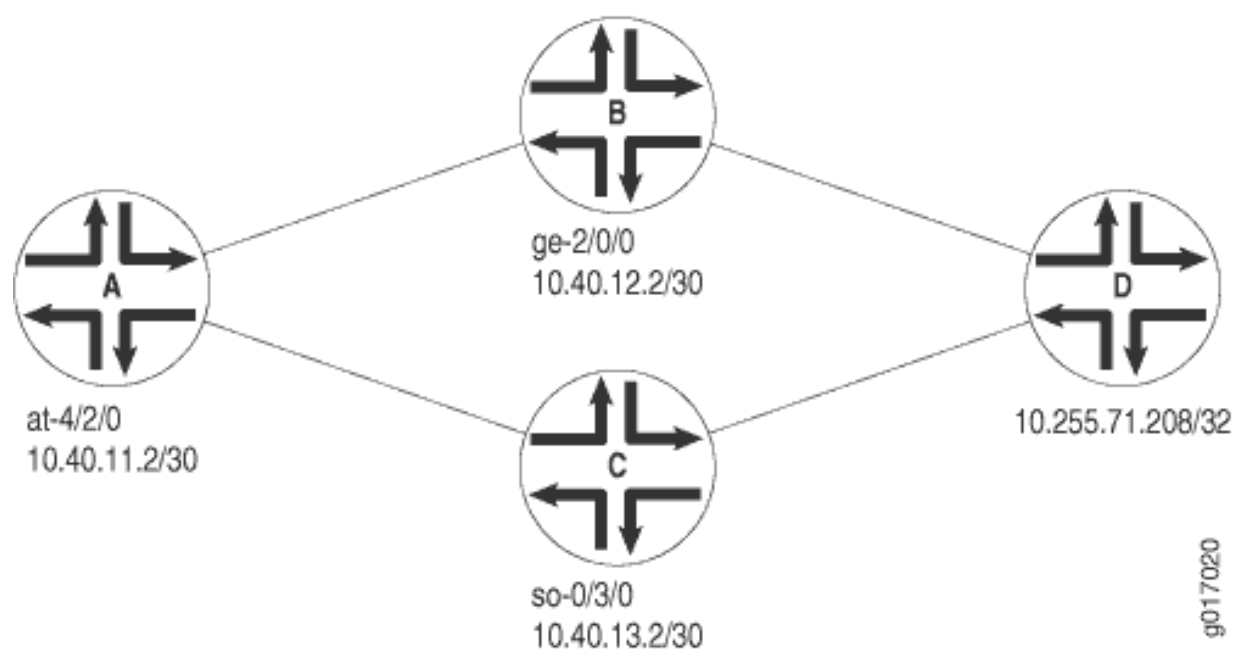
Forwarding Policy Options Overview

Example: Configuring CoS-Based Forwarding

Router A has two routes to destination 10.255.71.208 on Router D. One route goes through Router B, and the other goes through Router C, as shown in [Figure 5 on page 171](#).

Configure Router A with CoS-based forwarding (CBF) to select Router B for queue 0 and queue 2, and Router C for queue 1 and queue 3.

Figure 5: Sample CoS-Based Forwarding



When you configure CBF with OSPF as the IGP, you must specify the next hop as an interface name, not as an IPv4 or IPv6 address. The next hops in this example are specified as ge-2/0/0.0 and so-0/3/0.0.

```
[edit class-of-service]
forwarding-policy {
  next-hop-map my_cbf {
    forwarding-class be {
      next-hop ge-2/0/0.0;
    }
    forwarding-class ef {
      next-hop so-0/3/0.0;
    }
    forwarding-class af {
      next-hop ge-2/0/0.0;
    }
    forwarding-class nc {
      next-hop so-0/3/0.0;
    }
  }
}
classifiers {
  inet-precedence inet {
    forwarding-class be {
      loss-priority low code-points [ 000 100 ];
    }
    forwarding-class ef {
      loss-priority low code-points [ 001 101 ];
    }
    forwarding-class af {
      loss-priority low code-points [ 010 110 ];
    }
    forwarding-class nc {
      loss-priority low code-points [ 011 111 ];
    }
  }
}
forwarding-classes {
  queue 0 be;
  queue 1 ef;
  queue 2 af;
  queue 3 nc;
}
```

```

interfaces {
    at-4/2/0 {
        unit 0 {
            classifiers {
                inet-precedence inet;
            }
        }
    }
}

[edit policy-options]
policy-statement cbf {
    from {
        route-filter 10.255.71.208/32 exact;
    }
    then cos-next-hop-map my_cbf;
}

[edit routing-options]
graceful-restart;
forwarding-table {
    export cbf;
}

[edit interfaces]
traceoptions {
    file trace-intf size 5m world-readable;
    flag all;
}
so-0/3/0 {
    unit 0 {
        family inet {
            address 10.40.13.1/30;
        }
        family iso;
        family mpls;
    }
}
ge-2/0/0 {
    unit 0 {
        family inet {
            address 10.40.12.1/30;
        }
    }
}

```

```

        family iso;
        family mpls;
    }
}
at-4/2/0 {
    atm-options {
        vpi 1 {
            maximum-vcs 1200;
        }
    }
    unit 0 {
        vci 1.100;
        family inet {
            address 10.40.11.2/30;
        }
        family iso;
        family mpls;
    }
}

```

RELATED DOCUMENTATION

| *Forwarding Policy Options Overview*

Example: Configuring Forwarding Classes

IN THIS SECTION

- [Requirements | 175](#)
- [Overview | 175](#)
- [Example 1: Configuring Forwarding Classes for Switches Except QFX10000 | 177](#)
- [Example 2: Configuring Forwarding Classes for QFX10000 Switches | 179](#)

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 multidestination (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 multidestination 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 multidestination 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 multidestination forwarding class. You can define the remaining forwarding classes and configure them as unicast or multidestination forwarding classes by mapping them to unicast or multidestination queues. The type of queue, unicast or multidestination, 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 | 178](#)

Configuration

Step-by-Step Procedure

[Table 45 on page 177](#) shows the configuration forwarding-class-to-queue mapping for this example:

Table 45: 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 | 178](#)

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 | 180

Configuration

Step-by-Step Procedure

[Table 46 on page 179](#) shows the configuration forwarding-class-to-queue mapping for this example:

Table 46: Forwarding-Class-to-Queue Example Configuration on QFX10000

Forwarding Class	Queue
best-effort	0
be1	1
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 | 180](#)

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 \(ETS\) | 445](#)

[Defining CoS Forwarding Classes | 162](#)

[Monitoring CoS Forwarding Classes | 189](#)

Overview of CoS Changes Introduced in Junos OS Release 11.3

Overview of CoS Changes Introduced in Junos OS Release 12.2

[Understanding CoS Forwarding Classes | 155](#)

[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 multideestination 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 \(ETS\) | 438](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[Example: Configuring Forwarding Class Sets | 184](#)

[Defining CoS Forwarding Class Sets | 183](#)

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 \(ETS\) | 445](#)

[Example: Configuring Forwarding Class Sets | 184](#)

[Defining CoS Queue Schedulers | 346](#)

[Defining CoS Traffic Control Profiles \(Priority Group Scheduling\) | 412](#)

[Understanding CoS Forwarding Class Sets \(Priority Groups\) | 181](#)

Example: Configuring Forwarding Class Sets

IN THIS SECTION

- [Requirements | 185](#)
- [Overview | 185](#)
- [Verification | 187](#)

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 \(ETS\)" on page 445](#) provides a complete example of how to configure ETS, including forwarding class and scheduling configuration. [Table 47 on page 186](#) shows the configuration components for this example:

Table 47: 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
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 | 187](#)
- [Verifying the Egress Interface Configuration | 188](#)

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 \(ETS\) | 445](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Defining CoS Forwarding Class Sets | 183](#)

[Understanding CoS Forwarding Class Sets \(Priority Groups\) | 181](#)

Monitoring CoS Forwarding Classes

IN THIS SECTION

- Purpose | 189
- Action | 189
- Meaning | 189

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 48 on page 190](#) summarizes key output fields on switches that use different forwarding classes and output queues for unicast and multideestination traffic.

Table 48: 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"> • best-effort—Provides no special CoS handling of packets. Loss priority is typically not carried in a CoS value. • fcoe—Provides guaranteed delivery for Fibre Channel over Ethernet (FCoE) traffic. • no-loss—Provides guaranteed delivery for TCP lossless traffic • network-control—Packets can be delayed but not dropped. <p>By default, the following multidestination forwarding class is assigned to queue 8:</p> <ul style="list-style-type: none"> • mcast—Provides no special CoS handling of packets.
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"> • Queue 0—best-effort • Queue 3—fcoe • Queue 4—no-loss • Queue 7—network-control • Queue 8—mcast

Table 48: 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"> • Disabled—The forwarding class is configured for lossy transport (packets might drop during periods of congestion) • Enabled—The forwarding class is configured for lossless transport <p>NOTE: 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 49 on page 192 summarizes key output fields on switches that use the same forwarding classes and output queues for unicast and multidestination traffic.

Table 49: 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"> • best-effort—Provides no special CoS handling of packets. Loss priority is typically not carried in a CoS value. • fcoe—Provides guaranteed delivery for Fibre Channel over Ethernet (FCoE) traffic. • no-loss—Provides guaranteed delivery for TCP lossless traffic • network-control—Packets can be delayed but not dropped.
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"> • Queue 0—best-effort • Queue 3—fcoe • Queue 4—no-loss • Queue 7—network-control

Table 49: 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">• Disabled—The forwarding class is configured for lossy transport (packets might drop during periods of congestion).• Enabled—The forwarding class is configured for lossless transport. <p>NOTE: 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>

CHAPTER 7

Lossless Traffic Flows, Ethernet PAUSE Flow Control, and PFC

IN THIS CHAPTER

- Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 194
- Configuring CoS PFC (Congestion Notification Profiles) | 216
- Understanding CoS Flow Control (Ethernet PAUSE and PFC) | 220
- Enabling and Disabling CoS Symmetric Ethernet PAUSE Flow Control | 233
- Configuring CoS Asymmetric Ethernet PAUSE Flow Control | 234
- Understanding PFC Functionality Across Layer 3 Interfaces | 236
- Example: Configuring PFC Across Layer 3 Interfaces | 239
- Understanding PFC Using DSCP at Layer 3 for Untagged Traffic | 266
- Configuring DSCP-based PFC for Layer 3 Untagged Traffic | 268

Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows

IN THIS SECTION

- Default Lossless Priority Configuration | 195
- Configuring Lossless Priorities | 198
- Configuration Rules and Recommendations | 213
- Lossless Transport Features Introduced in Junos OS Release 12.3 (Legacy Non-ELS CLI) | 214
- Backward Compatibility with Junos OS Releases Earlier Than Release 12.3 (Legacy Non-ELS CLI) | 214

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 50 on page 197](#) summarizes the default forwarding classes and their mapping to output queues, IEEE 802.1p priorities, and drop attributes.

Table 50: 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](#)" on page 684 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 51 on page 199](#) 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 51: 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	<p>The fcoe and no-loss forwarding classes are lossless.</p> <p>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.</p>
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.

Table 51: FCoE and No-Loss Forwarding Class Configuration in Junos OS Release 12.3 (Continued)

Explicit (User-Configured) or Default Forwarding Class Configuration	Packet Drop Attribute	Result and Notes
Explicit, configured in Junos OS Release 12.2 or earlier	Not specified (packet drop attribute was not available before Junos OS Release 12.3)	<p>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.</p> <p>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.</p>

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 52 on page 205](#) 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 52: 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

Table 52: Default Output Flow Control Profile (Continued)

IEEE 802.1p Priority Specified in Received PFC Frame	Paused Output Queue
4 (100)	4
5 (101)	—
6 (110)	—
7 (111)	—

[Table 53 on page 206](#) 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 53: 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

Table 53: User-Configured Output Flow Control Profile (Continued)

IEEE 802.1p Priority Specified in Received PFC Frame	Paused Output Queue
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" on page 236](#) for a conceptual overview of how to enable PFC on traffic across Layer 3 interfaces. See ["Example: Configuring PFC Across Layer 3 Interfaces" on page 239](#) 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" on page 500](#) 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*.

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 54 on page 211](#) 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 54: 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 54: 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 54: 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 multidestination traffic, you cannot configure flow control to pause a multidestination 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 multidestination traffic, forwarding classes mapped to multidestination queues (queues 8 through 11) cannot have the no-loss packet drop attribute. (Multidestination 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*).

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](#)" on page 684 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 | 500](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

[Understanding PFC Functionality Across Layer 3 Interfaces | 236](#)

[Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic \(FCoE Transit Switch\) | 608](#)

[Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface | 620](#)

[Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces | 633](#)

[Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications \(FCoE and iSCSI\) | 652](#)

[Example: Configuring PFC Across Layer 3 Interfaces | 239](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

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 | 524](#)

[Assigning CoS Components to Interfaces | 87](#)

Monitoring Interfaces That Have CoS Components

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

[Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 194](#)

Understanding CoS Flow Control (Ethernet PAUSE and PFC)

IN THIS SECTION

- [General Information about Ethernet PAUSE and PFC and When to Use Them | 221](#)
- [Ethernet PAUSE | 222](#)
- [PFC | 227](#)
- [Lossless Transport Support Summary | 231](#)

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 55 on page 224](#) describes the configured flow control state when you set the receive (Rx) and transmit (Tx) buffers on an interface:

Table 55: 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.)
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 56 on page 224](#):

Table 56: 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.

Table 56: Flow Control State Advertised to the Connected Peer (Autonegotiation) (Continued)

Rx Buffer State	Tx Buffer State	PAUSE Bit	ASM_DIR Bit	Description
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).
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 57 on page 226. 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 57: 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

Table 57: 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	0	Disable Ethernet PAUSE transmit and receive	Disable Ethernet PAUSE transmit and receive
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 57 on page 226](#) 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 58 on page 228](#):

Table 58: Default PFC Priority to Queue and Forwarding Class Mapping

IEEE 802.1p Priority (Code Point)	Queue	Forwarding Class
0 (000)	0	best-effort

Table 58: Default PFC Priority to Queue and Forwarding Class Mapping (*Continued*)

IEEE 802.1p Priority (Code Point)	Queue	Forwarding Class
1 (001)	1	best-effort
2 (010)	2	best-effort
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 multdestination traffic, you can map a unicast queue (queue 0 through 7) and a multdestination 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 | 482](#)

[Understanding CoS Explicit Congestion Notification](#)

[Configuring CoS PFC \(Congestion Notification Profiles\)](#)

[Example: Configuring CoS PFC for FCoE Traffic | 524](#)

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 | 234](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

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 | 233](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

Understanding PFC Functionality Across Layer 3 Interfaces

Priority-based flow control (PFC) allows you to select traffic flows within a link and pause them, so that the output queues associated with the flows do not overflow and drop packets. (PFC is more granular than Ethernet PAUSE, which pauses all traffic on a physical link.) PFC helps you configure lossless transport for traffic flows across a data center bridging network.

However, you might want to create a traffic flow that losslessly traverses the Layer 2 data center bridging network *and* also losslessly traverses a Layer 3 network that connects Ethernet hosts in different Layer 2 networks. On a QFX5210, QFX5200, QFX5110, QFX5100, EX4600, or QFX10000 switch running the Enhanced Layer 2 Software (ELS) CLI, in addition to configuring PFC on Layer 2 (bridging) interfaces, you can configure PFC on VLAN-tagged traffic that traverses Layer 3 interfaces. This enables you to preserve the lossless characteristics that PFC provides on VLAN-tagged traffic, even when the traffic crosses Layer 3 interfaces that connect two Layer 2 networks.

NOTE: This topic is applicable for VLAN-tagged traffic only. Starting in Junos OS Release 17.4R1, QFX5110, QFX5200, and QFX5210 switches also support DSCP-based PFC for *untagged* traffic on Layer 3 interfaces and Layer 2 access interfaces. DSCP-based PFC uses a DSCP classifier to classify the traffic based on a 6-bit DSCP value that is mapped to a 3-bit PFC priority value. For details on using DSCP-based PFC on supporting switches, see *Understanding PFC Using DSCP at Layer 3 for Untagged Traffic*.



Video: [Preserving Lossless Behavior on an SDN or Overlay Network](#)

PFC works the same way across Layer 3 interfaces as it works across Layer 2 interfaces. When an output queue buffer reaches a certain fill level threshold, the switch sends a PFC pause message to the

connected peer to pause transmission of the traffic on which PFC is enabled. Pausing the incoming traffic prevents the queue buffer from overflowing and dropping packets, just as on Layer 2 interfaces. When the queue buffer fill level decreases below a certain threshold, the interface sends a message to the connected peer to restart traffic transmission.

Although PFC is a data center bridging technology, PFC also works on Layer 3 interfaces because PFC operates at the queue level. When you use an IEEE 802.1p classifier to classify incoming traffic (map incoming traffic to a forwarding class and a loss priority based on the IEEE 802.1p code point in the Ethernet frame header) and you enable PFC on the appropriate priority (IEEE 802.1p code point), PFC works on Layer 2 and Layer 3 interfaces.

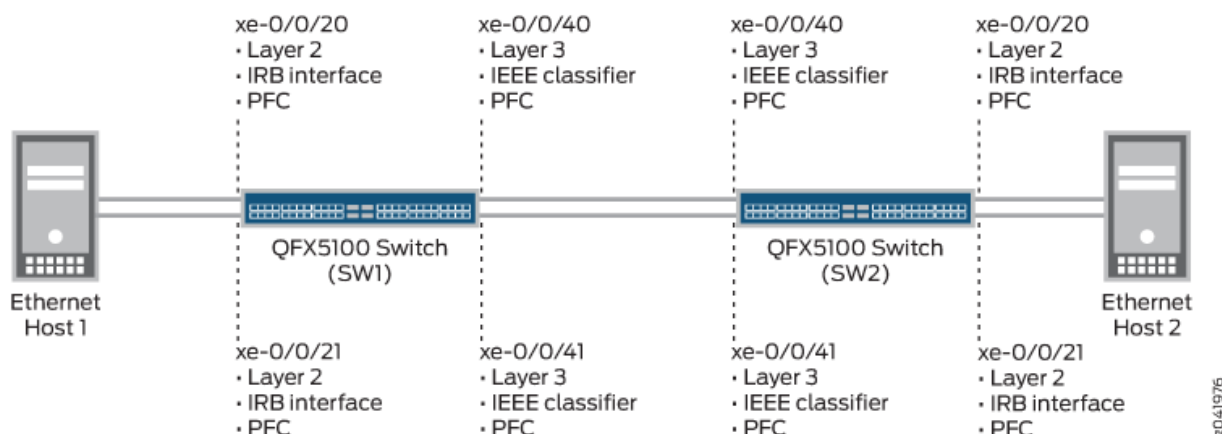
NOTE: Lossless VLAN-tagged traffic on Layer 3 interfaces *must* use an IEEE 802.1p classifier to classify incoming traffic, because PFC does not use DSCP or DSCP IPv6 code points to identify VLAN-tagged traffic for flow control. PFC cannot pause traffic flows unless the incoming traffic is classified by an IEEE 802.1p classifier. Do not apply a DSCP (or a DSCP IPv6) classifier to Layer 3 VLAN-tagged traffic on which you want to enable PFC.

Because PFC functionality relies on the mapping (classifying) of incoming traffic to IEEE 802.1p code points and on enabling PFC on the correct code point(s) at each interface, you must ensure that incoming traffic has the correct 3-bit IEEE 802.1p code point (priority) in the priority code point (PCP) field of the Ethernet frame header (sometimes known as the CoS bits).

NOTE: Layer 3 interfaces do not support FCoE traffic. FCoE traffic must use Layer 2 interfaces and cannot use Layer 3 interfaces. Therefore, you cannot enable PFC on FCoE traffic across Layer 3 interfaces.

[Figure 6 on page 238](#) shows a topology in which two Ethernet hosts in Layer 2 networks communicate across a Layer 3 network, with PFC enabled on all of the Layer 2 and Layer 3 switch interfaces.

Figure 6: Enabling PFC Across Layer 3 Interface Hops



The Ethernet host-facing interfaces (xe-0/0/20 and xe-0/0/21 on both switches) and the Layer 3 network-facing interfaces (interfaces xe-0/0/40 and xe-0/0/41 on both switches) require different interface configurations to enable PFC on the Layer 3 interfaces. In addition, the class of service (CoS) for each interface must be configured correctly, including enabling PFC on the traffic that you want to treat as lossless traffic:

Ethernet-host facing interfaces (xe-0/0/20 and xe-0/0/21) require the following configuration:

- Set interfaces as family ethernet-switching
- Set the interface mode as trunk mode
- Create VLANs to carry the traffic
- Create IRB interfaces to place the Layer 2 VLAN traffic on Layer 3 for transport between IP networks
- Create an IEEE 802.1p classifier to classify incoming traffic into the correct forwarding class, based on the IEEE 802.1p code point
- Create a congestion notification profile (CNP) to configure PFC on the IEEE 802.1p code point of the traffic that you want treat as lossless traffic
- Apply the classifier and the CNP to the Layer 2 interfaces
- Configure CoS: lossless forwarding classes, hierarchical port scheduling (also known as enhanced transmission selection), or direct port scheduling, depending on your switch, and apply it to the Layer 2 interfaces

Layer 3 IP network-facing interfaces (xe-0/0/40 and xe-0/0/41) require the following configuration:

- Set interfaces as family inet
- Set VLAN tagging on the interfaces

- Create VLANs to carry the traffic
- Create an IEEE 802.1p classifier to classify incoming traffic into the correct forwarding class, based on the IEEE 802.1p code point (do not use a DSCP or DSCP IPv6 classifier)
- Create a congestion notification profile (CNP) to configure PFC on the IEEE 802.1p code point of the traffic that you want treat as lossless traffic on the Layer 3 interfaces
- Apply the IEEE 802.1p classifier and the CNP to the Layer 3 interfaces
- Configure CoS: lossless forwarding classes, hierarchical port scheduling (enhanced transmission selection), or direct port scheduling, depending on your switch, and apply it to the Layer 3 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.

When you configure the Layer 2 and Layer 3 interfaces correctly, the switch enables PFC on the traffic between Ethernet Host 1 and Ethernet Host 2 across the entire path between the two hosts. If any output queue in the path on which PFC is enabled experiences congestion, PFC pauses the traffic and prevents packet loss for the flow.

RELATED DOCUMENTATION

[Example: Configuring PFC Across Layer 3 Interfaces | 239](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

[Understanding Integrated Routing and Bridging](#)

Example: Configuring PFC Across Layer 3 Interfaces

IN THIS SECTION

- [Requirements | 240](#)
- [Overview | 240](#)
- [Configuration | 246](#)

Priority-based flow control (PFC) helps ensure lossless transport across data center bridging interfaces by pausing incoming traffic when output queue buffers fill to a certain threshold. On a QFX5210, QFX5200, QFX5110, QFX5100, EX4600, or QFX10000 switch running the Enhanced Layer 2 Software (ELS) CLI, in addition to configuring PFC on Layer 2 (bridging) interfaces, you can configure PFC on VLAN-tagged traffic that traverses Layer 3 interfaces. This enables you to preserve the lossless characteristics that PFC provides on VLAN-tagged traffic, even when the traffic crosses Layer 3 interfaces that connect two Layer 2 networks.

NOTE: This topic is applicable for VLAN-tagged traffic only. Starting in Junos OS Release 17.4R1, QFX5110, QFX5200, and QFX5210 switches also support DSCP-based PFC for *untagged* traffic on Layer 3 interfaces and Layer 2 access interfaces. DSCP-based PFC uses a DSCP classifier to classify the traffic based on a 6-bit DSCP value that is mapped to a 3-bit PFC priority value. For details on configuring DSCP-based PFC on supporting switches, see *Configuring DSCP-based PFC for Layer 3 Untagged Traffic*.

Requirements

This example uses the following hardware and software components:

- Two switches
- Junos OS Release 13.2 or later for the QFX Series
- Two Ethernet hosts

Overview

IN THIS SECTION

- Topology | 241

On a network that uses two QFX5210, QFX5200, QFX5110, QFX5100, EX4600, or QFX10000 switches to connect hosts on two different Ethernet networks across a Layer 3 network, to configure PFC across the Layer 2 and Layer 3 interfaces, you must:

- Configure the Layer 2 and Layer 3 interfaces on the switches
- Configure VLANs to carry the traffic across the Layer 2 and Layer 3 networks
- Configure integrated routing and bridging (IRB) interfaces on the Layer 2 interfaces to move the Layer 2 VLAN traffic to Layer 3
- Configure and apply the appropriate classifiers to the interfaces
- Configure and apply congestion notification profiles (CNPs) on the interfaces to enable PFC on the traffic that you want to be lossless

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 lossless forwarding classes and either hierarchical port scheduling (also known as enhanced transmission selection) or direct port scheduling, depending on your switch, on the interfaces

NOTE: PFC operates at the queue level, based on the IEEE 802.1p code point in the priority code point (PCP) field of the Ethernet frame header (sometimes known as the CoS bits). For this reason, VLAN-tagged traffic on Layer 3 interfaces on which you want to enable PFC must use an IEEE 802.1p classifier to map incoming traffic to forwarding classes (which are in turn mapped to output queues) and loss priorities. You cannot use a DSCP or DSCP IPv6 classifier to classify Layer 3 traffic if you want to enable PFC on VLAN-tagged traffic flows.

Topology

Figure 7 on page 242 shows the topology for this example.

Figure 7: Enabling PFC Across Layer 3 Interface Hops

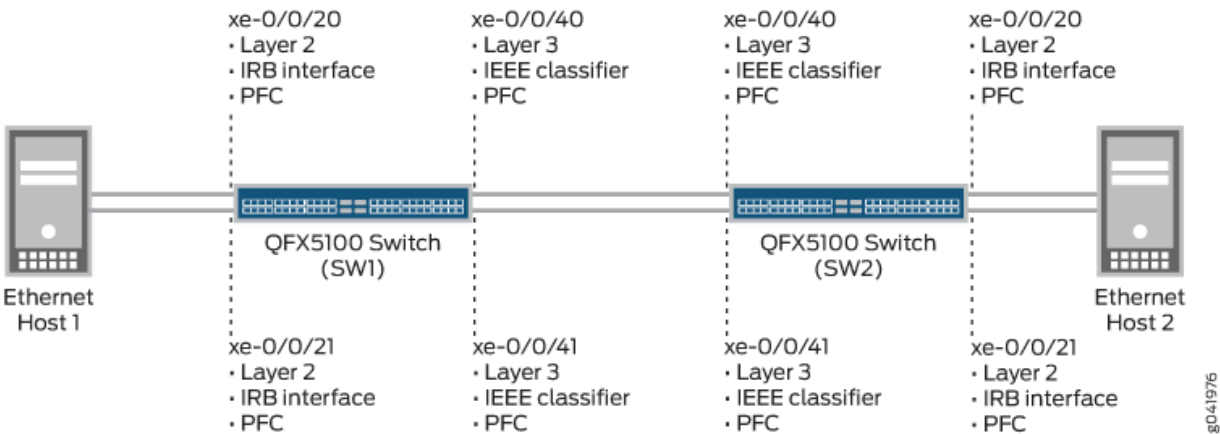


Table 59 on page 242 shows the configuration components for this example. On the two switches, the Ethernet host-facing interfaces use the same interface names and configuration, and the Layer 3 network-facing interfaces use the same interface names and configuration.

Table 59: Components of the PFC Across Layer 3 Interfaces Topology

Component	Settings
Hardware	Two switches, Switch SW1 and Switch SW2. Two Ethernet hosts
Layer 3 interfaces (xe-0/0/40 and xe-0/0/41) and VLANs	Interface xe-0/0/40: <ul style="list-style-type: none">Interface family—inetInterface IP address—100.103.1.2/24VLAN tagging—enabledInterface VLAN ID—103 Interface xe-0/0/41: <ul style="list-style-type: none">Interface family—inetInterface IP address—100.104.1.2/24VLAN tagging—enabledInterface VLAN ID—104

Table 59: Components of the PFC Across Layer 3 Interfaces Topology (*Continued*)

Component	Settings
Layer 2 interfaces (xe-0/0/20 and xe-0/0/21) and VLAN membership	Family: Ethernet switching Interface mode—trunk Interface xe-0/0/20 VLAN membership—vlan105 Interface xe-0/0/21 VLAN membership—vlan106
VLANs for the IRB interfaces	VLAN unit 105—family inet, IP address 100.105.1.1/24 VLAN unit 106—family inet, IP address 100.106.1.1/24
Layer 2 IRB interfaces	Interface xe-0/0/20: <ul style="list-style-type: none"> • IRB interface unit—105 • IRB interface family—inet • IRB interface IP address—100.105.1.1/24 • IRB interface VLAN ID—105 • Layer 3 interface name—irb.105 Interface xe-0/0/21: <ul style="list-style-type: none"> • IRB interface unit—106 • IRB interface family—inet • IRB interface IP address—100.106.1.1/24 • IRB interface VLAN ID—106 • Layer 3 interface name—irb.106

Table 59: Components of the PFC Across Layer 3 Interfaces Topology (Continued)

Component	Settings
Forwarding classes (both switches)	<p>Name—lossless-3 Queue mapping—queue 3 Packet drop attribute—no-loss</p> <p>Name—lossless-4 Queue mapping—queue 4 Packet drop attribute—no-loss</p> <p>NOTE: Matching the forwarding class names (lossless-3 and lossless-4) to the queue number and to the classified IEEE 802.1p code point (priority) creates a configuration that is logical and easy to map because the forwarding class, queue, and priority all use the same number.</p> <p>Name—all-others Queue mapping—queue 0 Packet drop attribute—none</p> <p>NOTE: The forwarding class <i>all-others</i> is for best-effort traffic that traverses the interfaces.</p>
Layer 2 interface behavior aggregate (BA) classifier	<p>Name—lossless-3-4-ieee Forwarding class lossless-3—mapped to code point 011 (IEEE 802.1p priority 3) and a packet loss priority of low Forwarding class lossless-4—mapped to code point 100 (IEEE 802.1p priority 4) and a packet loss priority of low</p> <p>Apply the Layer 2 IEEE 802.1p classifier to both the Layer 2 and the Layer 3 interfaces (xe-0/0/20, xe-0/0/21, xe-0/0/40, and xe-0/0/41).</p>
Congestion notification profile (PFC, both switches)	<p>Name—lossless-cnp PFC enabled on IEEE 802.1p code points—011 (lossless-3 forwarding class and priority), 100 (lossless-4 forwarding class and priority)</p> <p>Apply the CNP to both the Layer 2 and the Layer 3 interfaces (xe-0/0/20, xe-0/0/21, xe-0/0/40, and xe-0/0/41) to enable PFC on IEEE 802.1p code points 011 and 100.</p>

Table 59: Components of the PFC Across Layer 3 Interfaces Topology (*Continued*)

Component	Settings
Enhanced transmission selection (ETS) hierarchical port scheduling (only if using ETS)	<p>Hierarchical port scheduling (ETS) includes configuring:</p> <ul style="list-style-type: none"> • Schedulers to assign bandwidth to traffic • Scheduler mapping to forwarding classes • Grouping of the forwarding classes (priorities) in forwarding class sets (priority groups) • A traffic control profile to assign bandwidth to the forwarding class set and to associate the forwarding class set with the scheduler mapping <p>Hierarchical port scheduling also includes applying the hierarchical scheduler (defined in the traffic control profile) to the interfaces.</p> <p>This example focuses on configuring PFC across the Layer 2 and Layer 3 interfaces. To maintain this focus, this example includes the CLI statements needed to configure hierarchical port scheduling, but does not include descriptive explanations of the configuration. The <i>Related Documentation</i> section provides links to example documents that show how to configure hierarchical port scheduling.</p> <p>Apply the scheduling configuration to both the Layer 2 and the Layer 3 interfaces (xe-0/0/20, xe-0/0/21, xe-0/0/40, and xe-0/0/41).</p>
Direct port scheduling (only if using port scheduling instead of ETS)	<p>Direct port scheduling includes configuring:</p> <ul style="list-style-type: none"> • Schedulers to assign bandwidth to traffic • Scheduler mapping to forwarding classes <p>Port scheduling also includes applying the scheduler map to the interfaces.</p> <p>This example focuses on configuring PFC across the Layer 2 and Layer 3 interfaces. To maintain this focus, this example includes the CLI statements needed to configure direct port scheduling, but does not include descriptive explanations of the configuration. The <i>Related Documentation</i> section provides links to example documents that show how to configure port scheduling.</p> <p>Apply the scheduling configuration to both the Layer 2 and the Layer 3 interfaces (xe-0/0/20, xe-0/0/21, xe-0/0/40, and xe-0/0/41).</p>

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 246](#)
- [Common Configuration \(Applies to ETS Hierarchical Scheduling and to Port Scheduling\) | 248](#)
- [ETS Hierarchical Scheduling Configuration | 251](#)
- [Port Scheduling Configuration | 252](#)
- [Results | 252](#)

CLI Quick Configuration

To configure PFC across Layer 3 interfaces, copy the following commands, paste them in a text file, remove the 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 same configuration applies to both Switch SW1 and Switch SW2. 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 (Applies to ETS Hierarchical Scheduling and to Port Scheduling)

```
set interfaces xe-0/0/40 vlan-tagging
set interfaces xe-0/0/40 unit 0 vlan-id 103
set interfaces xe-0/0/40 unit 0 family inet address 100.103.1.2/24
set interfaces xe-0/0/41 vlan-tagging
set interfaces xe-0/0/41 unit 0 vlan-id 104
set interfaces xe-0/0/41 unit 0 family inet address 100.104.1.2/24
set interfaces xe-0/0/20 unit 0 family ethernet-switching interface-mode trunk
set interfaces xe-0/0/20 unit 0 family ethernet-switching vlan members vlan105
set interfaces xe-0/0/21 unit 0 family ethernet-switching interface-mode trunk
set interfaces xe-0/0/21 unit 0 family ethernet-switching vlan members vlan106
set interfaces irb unit 105 family inet address 100.105.1.1/24
set interfaces irb unit 106 family inet address 100.106.1.1/24
set vlans vlan105 vlan-id 105
set vlans vlan106 vlan-id 106
set vlans vlan105 l3-interface irb.105
set vlans vlan106 l3-interface irb.106
set class-of-service forwarding-classes class lossless-3 queue-num 3 no-loss
```

```

set class-of-service forwarding-classes class lossless-4 queue-num 4 no-loss
set class-of-service forwarding-classes class all-others queue-num 0
set class-of-service classifiers ieee-802.1 lossless-3-4-ieee forwarding-class lossless-3 loss-
priority low code-points 011
set class-of-service classifiers ieee-802.1 lossless-3-4-ieee forwarding-class lossless-4 loss-
priority low code-points 100
set class-of-service congestion-notification-profile lossless-cnp input ieee-802.1 code-point
011 pfc
set class-of-service congestion-notification-profile lossless-cnp input ieee-802.1 code-point
100 pfc
set class-of-service schedulers lossless_sch transmit-rate 6g
set class-of-service schedulers lossless_sch shaping-rate percent 100
set class-of-service schedulers all-others_sch transmit-rate 4g
set class-of-service scheduler-maps lossless_map forwarding-class lossless-3 scheduler
lossless_sch
set class-of-service scheduler-maps lossless_map forwarding-class lossless-4 scheduler
lossless_sch
set class-of-service scheduler-maps all-others_map forwarding-class all-others scheduler all-
others_sch
set class-of-service interfaces xe-0/0/20 congestion-notification-profile lossless-cnp
set class-of-service interfaces xe-0/0/20 unit 0 classifiers ieee-802.1 lossless-3-4-ieee
set class-of-service interfaces xe-0/0/21 congestion-notification-profile lossless-cnp
set class-of-service interfaces xe-0/0/21 unit 0 classifiers ieee-802.1 lossless-3-4-ieee
set class-of-service interfaces xe-0/0/40 congestion-notification-profile lossless-cnp
set class-of-service interfaces xe-0/0/40 classifiers ieee-802.1 lossless-3-4-ieee
set class-of-service interfaces xe-0/0/41 congestion-notification-profile lossless-cnp
set class-of-service interfaces xe-0/0/41 classifiers ieee-802.1 lossless-3-4-ieee

```

Configuration for ETS Hierarchical Scheduling

The ETS-specific portion of this example configures forwarding class set (priority group) membership and priority group CoS settings (traffic control profile), and assigns the priority group and its CoS configuration to the interfaces.

```

set class-of-service forwarding-class-sets lossless_fc_set class lossless-3
set class-of-service forwarding-class-sets lossless_fc_set class lossless-4
set class-of-service forwarding-class-sets all-others_fc_set class all-others
set class-of-service traffic-control-profiles lossless_tcp scheduler-map lossless_map
set class-of-service traffic-control-profiles lossless_tcp guaranteed-rate percent 60
set class-of-service traffic-control-profiles lossless_tcp shaping-rate percent 100
set class-of-service traffic-control-profiles all-others_tcp scheduler-map all-others_map
set class-of-service traffic-control-profiles all-others_tcp guaranteed-rate percent 40

```

```

set class-of-service interfaces xe-0/0/20 forwarding-class-set lossless_fc_set output-traffic-
control-profile lossless_tcp
set class-of-service interfaces xe-0/0/20 forwarding-class-set all-others_fc_set output-traffic-
control-profile all-others_tcp
set class-of-service interfaces xe-0/0/21 forwarding-class-set lossless_fc_set output-traffic-
control-profile lossless_tcp
set class-of-service interfaces xe-0/0/21 forwarding-class-set all-others_fc_set output-traffic-
control-profile all-others_tcp
set class-of-service interfaces xe-0/0/40 forwarding-class-set lossless_fc_set output-traffic-
control-profile lossless_tcp
set class-of-service interfaces xe-0/0/40 forwarding-class-set all-others_fc_set output-traffic-
control-profile all-others_tcp
set class-of-service interfaces xe-0/0/41 forwarding-class-set lossless_fc_set output-traffic-
control-profile lossless_tcp
set class-of-service interfaces xe-0/0/41 forwarding-class-set all-others_fc_set output-traffic-
control-profile all-others_tcp

```

Configuration for Port Scheduling

The port-scheduling-specific portion of this example assigns the scheduler maps (which set the CoS treatment of the forwarding classes in the scheduler map) to the interfaces.

```

[edit class-of-service]
set interfaces xe-0/0/20 scheduler-map lossless_map
set interfaces xe-0/0/20 scheduler-map all-others_map
set interfaces xe-0/0/21 scheduler-map lossless_map
set interfaces xe-0/0/21 scheduler-map all-others_map
set interfaces xe-0/0/40 scheduler-map lossless_map
set interfaces xe-0/0/40 scheduler-map all-others_map
set interfaces xe-0/0/41 scheduler-map lossless_map
set interfaces xe-0/0/41 scheduler-map all-others_map

```

Common Configuration (Applies to ETS Hierarchical Scheduling and to Port Scheduling)

Step-by-Step Procedure

The following step-by-step procedure shows you how to configure the VLANs, IRB interfaces, lossless forwarding classes, classifiers, PFC settings to enable PFC across Layer 3 interfaces, and the queue scheduling configuration common to ETS and direct port scheduling. For completeness, the ETS hierarchical port scheduling and direct port scheduling configurations are included separately, in the

following procedures, but without explanatory text. See the *Related Documentation* links for detailed examples of the scheduling elements of the configuration.

1. Configure the Layer 3 interface VLANs and IP addresses:

```
[edit interfaces]
user@switch# set xe-0/0/40 vlan-tagging
user@switch# set xe-0/0/40 unit 0 vlan-id 103
user@switch# set xe-0/0/40 unit 0 family inet address 100.103.1.2/24
user@switch# set xe-0/0/41 vlan-tagging
user@switch# set xe-0/0/41 unit 0 vlan-id 104
user@switch# set xe-0/0/41 unit 0 family inet address 100.104.1.2/24
```

2. Configure the Layer 2 interface VLAN membership and interface mode:

```
[edit interfaces]
user@switch# set xe-0/0/20 unit 0 family ethernet-switching interface-mode trunk
user@switch# set xe-0/0/20 unit 0 family ethernet-switching vlan members vlan105
user@switch# set xe-0/0/21 unit 0 family ethernet-switching interface-mode trunk
user@switch# set xe-0/0/21 unit 0 family ethernet-switching vlan members vlan106
```

3. Configure the IRB interfaces and VLANs to transport incoming Layer 2 traffic assigned to VLANs vlan105 (of which interface xe-0/0/20 is a member) and vlan106 (of which interface xe-0/0/21 is a member) across Layer 3:

```
[edit]
user@switch# set interfaces irb unit 105 family inet address 100.105.1.1/24
user@switch# set interfaces irb unit 106 family inet address 100.106.1.1/24
user@switch# set vlans vlan105 vlan-id 105
user@switch# set vlans vlan106 vlan-id 106
user@switch# set vlans vlan105 l3-interface irb.105
user@switch# set vlans vlan106 l3-interface irb.106
```

4. Configure the lossless forwarding classes and a best-effort forwarding class for any other traffic that might use the interfaces:

```
[edit class-of-service]
user@switch# set forwarding-classes class lossless-3 queue-num 3 no-loss
```

```

user@switch# set forwarding-classes class lossless-4 queue-num 4 no-loss
user@switch# set forwarding-classes class all-others queue-num 0

```

5. Configure the IEEE classifier for the Layer 2 and Layer 3 interfaces to classify incoming traffic into the lossless forwarding classes based on the IEEE 802.1p code point of the traffic:

```

[edit class-of-service classifiers]
user@switch# set ieee-802.1 lossless-3-4-ieee forwarding-class lossless-3 loss-priority low
code-points 011
user@switch# set ieee-802.1 lossless-3-4-ieee forwarding-class lossless-4 loss-priority low
code-points 100

```

6. Configure the CNP to enable PFC on the lossless priorities (the lossless forwarding classes mapped to IEEE 802.1p code points 3 and 4):

```

[edit class-of-service congestion-notification-profile]
user@switch# set lossless-cnp input ieee-802.1 code-point 011 pfc
user@switch# set lossless-cnp input ieee-802.1 code-point 100 pfc

```

7. Apply the Layer 2 IEEE 802.1p classifier and the CNP to the Layer 3 interfaces:

```

[edit class-of-service interfaces]
user@switch# set xe-0/0/40 classifiers ieee-802.1 lossless-3-4-ieee
user@switch# set xe-0/0/40 congestion-notification-profile lossless-cnp
user@switch# set xe-0/0/41 classifiers ieee-802.1 lossless-3-4-ieee
user@switch# set xe-0/0/41 congestion-notification-profile lossless-cnp

```

8. Apply the Layer 2 IEEE 802.1p classifier and the CNP to the Layer 2 interfaces:

```

[edit class-of-service interfaces]
user@switch# xe-0/0/20 unit 0 classifiers ieee-802.1 lossless-3-4-ieee
user@switch# xe-0/0/20 congestion-notification-profile lossless-cnp
user@switch# xe-0/0/21 unit 0 classifiers ieee-802.1 lossless-3-4-ieee
user@switch# xe-0/0/21 congestion-notification-profile lossless-cnp

```

9. Configure queue scheduling to support the lossless configuration and map the schedulers to the forwarding classes (statements included here for completeness; see the *Related Documentation* links for detailed examples of scheduling configuration):

```
[edit class-of-service]
user@switch# set schedulers lossless_sch transmit-rate 6g
user@switch# set schedulers lossless_sch shaping-rate percent 100
user@switch# set schedulers all-others_sch transmit-rate 4g
user@switch# set scheduler-maps lossless_map forwarding-class lossless-3 scheduler
lossless_sch
user@switch# set scheduler-maps lossless_map forwarding-class lossless-4 scheduler
lossless_sch
user@switch# set scheduler-maps all-others_map forwarding-class all-others scheduler all-
others_sch
```

ETS Hierarchical Scheduling Configuration

Step-by-Step Procedure

1. Configure hierarchical scheduling to support the lossless configuration (included here for completeness; see the *Related Documentation* links for detailed examples of scheduling configuration) and apply it to the Layer 2 and Layer 3 interfaces:

```
[edit class-of-service interfaces]
user@switch# set forwarding-class-sets lossless_fc_set class lossless-3
user@switch# set forwarding-class-sets lossless_fc_set class lossless-4
user@switch# set forwarding-class-sets all-others_fc_set class all-others
user@switch# set traffic-control-profiles lossless_tcp scheduler-map lossless_map
user@switch# set traffic-control-profiles lossless_tcp guaranteed-rate percent 60
user@switch# set traffic-control-profiles lossless_tcp shaping-rate percent 100
user@switch# set traffic-control-profiles all-others_tcp scheduler-map all-others_map
user@switch# set traffic-control-profiles all-others_tcp guaranteed-rate percent 40
user@switch# set interfaces xe-0/0/20 forwarding-class-set lossless_fc_set output-traffic-
control-profile lossless_tcp
user@switch# set interfaces xe-0/0/20 forwarding-class-set all-others_fc_set output-traffic-
control-profile all-others_tcp
user@switch# set interfaces xe-0/0/21 forwarding-class-set lossless_fc_set output-traffic-
control-profile lossless_tcp
user@switch# set interfaces xe-0/0/21 forwarding-class-set all-others_fc_set output-traffic-
control-profile all-others_tcp
```



```

user@switch# set interfaces xe-0/0/40 forwarding-class-set lossless_fc_set output-traffic-
control-profile lossless_tcp
user@switch# set interfaces xe-0/0/40 forwarding-class-set all-others_fc_set output-traffic-
control-profile all-others_tcp
user@switch# set interfaces xe-0/0/41 forwarding-class-set lossless_fc_set output-traffic-
control-profile lossless_tcp
user@switch# set interfaces xe-0/0/41 forwarding-class-set all-others_fc_set output-traffic-
control-profile all-others_tcp

```

Port Scheduling Configuration

Step-by-Step Procedure

1. Apply port scheduling to support the lossless configuration on interfaces:

```

[edit class-of-service]
user@switch# set interfaces xe-0/0/20 scheduler-map lossless_map
user@switch# set interfaces xe-0/0/20 scheduler-map all-others_map
user@switch# set interfaces xe-0/0/21 scheduler-map lossless_map
user@switch# set interfaces xe-0/0/21 scheduler-map all-others_map
user@switch# set interfaces xe-0/0/40 scheduler-map lossless_map
user@switch# set interfaces xe-0/0/40 scheduler-map all-others_map
user@switch# set interfaces xe-0/0/41 scheduler-map lossless_map
user@switch# set interfaces xe-0/0/41 scheduler-map all-others_map

```

Results

Display the results of the interface, VLAN, and class-of-service configurations (the system shows only the explicitly configured parameters; it does not show default parameters). The results are valid for both Switch SW1 and Switch SW2 because the same configuration is used on both switches. The results are from the ETS hierarchical scheduling configuration, which show the more complex configuration. Direct port scheduling results would not show the traffic control profile or forwarding class set portions of the configuration, but 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). Other than that, the results are the same.

Display the results of the interface configuration:

```

user@switch# show configuration interfaces
xe-0/0/20 {
  unit 0 {
    family ethernet-switching {
      interface-mode trunk;
      vlan {
        members vlan105;
      }
    }
  }
}
xe-0/0/21 {
  unit 0 {
    family ethernet-switching {
      interface-mode trunk;
      vlan {
        members vlan106;
      }
    }
  }
}
xe-0/0/40 {
  vlan-tagging;
  unit 0 {
    vlan-id 103;
    family inet {
      address 100.103.1.2/24;
    }
  }
}
xe-0/0/41 {
  vlan-tagging;
  unit 0 {
    vlan-id 104;
    family inet {
      address 100.104.1.2/24;
    }
  }
}
irb {

```

```

    unit 105 {
        family inet {
            address 100.105.1.1/24;
        }
    }
    unit 106 {
        family inet {
            address 100.106.1.1/24;
        }
    }
}
vlan {
    unit 105 {
        family inet {
            address 100.105.1.1/24;
        }
    }
    unit 106 {
        family inet {
            address 100.106.1.1/24;
        }
    }
}
}

```

Display the results of the vlan configuration:

```

user@switch# show configuration vlans
vlan105 {
    vlan-id 105;
    l3-interface irb.105;
}
vlan106 {
    vlan-id 106;
    l3-interface irb.106;
}

```

Display the results of the class-of-service configuration:

```

user@switch# show configuration class-of-service
classifiers {

```

```

ieee-802.1 lossless-3-4-ieee {
    forwarding-class lossless-3 {
        loss-priority low code-points 011;
    }
    forwarding-class lossless-4 {
        loss-priority low code-points 100;
    }
}
}
forwarding-classes {
    class lossless-3 queue-num 3 no-loss;
    class lossless-4 queue-num 4 no-loss;
    class all-others queue-num 0;
}
traffic-control-profiles {
    lossless_tcp {
        scheduler-map lossless_map;
        shaping-rate percent 100;
        guaranteed-rate percent 60;
    }
    all-others_tcp {
        scheduler-map all-others_map;
        guaranteed-rate percent 40;
    }
}
forwarding-class-sets {
    lossless_fc_set {
        class lossless-3;
        class lossless-4;
    }
    all-others_fc_set {
        class all-others;
    }
}
congestion-notification-profile {
    lossless-cnp {
        input {
            ieee-802.1 {
                code-point 011 {
                    pfc;
                }
                code-point 100 {
                    pfc;
                }
            }
        }
    }
}

```



```

        all-others_fc_set {
            output-traffic-control-profile all-others_tcp;
        }
    }
    congestion-notification-profile lossless-cnp;
    classifiers {
        ieee-802.1 lossless-3-4-ieee;
    }
}
xe-0/0/41 {
    forwarding-class-set {
        lossless_fc_set {
            output-traffic-control-profile lossless_tcp;
        }
        all-others_fc_set {
            output-traffic-control-profile all-others_tcp;
        }
    }
    congestion-notification-profile lossless-cnp;
    classifiers {
        ieee-802.1 lossless-3-4-ieee;
    }
}
}
scheduler-maps {
    lossless_map {
        forwarding-class lossless-3 scheduler lossless_sch;
        forwarding-class lossless-4 scheduler lossless_sch;
    }
    all-others_map {
        forwarding-class all-others scheduler all-others_sch;
    }
}
}
schedulers {
    lossless_sch {
        transmit-rate 6g;
        shaping-rate percent 100;
    }
    all-others_sch {
        transmit-rate 4g;
    }
}
}

```

TIP: To quickly configure the switch, issue the `load merge terminal` command, and then copy the hierarchies and paste them into the switch terminal window.

Verification

IN THIS SECTION

- [Verifying the Interface Configuration | 258](#)
- [Verifying the VLAN Configuration | 260](#)
- [Verifying the PFC Configuration \(Congestion Notification Profile\) | 261](#)
- [Verify the Forwarding Class Configuration | 262](#)
- [Verifying the Classifier Configuration | 263](#)
- [Verifying the Interface CoS Configuration \(Hierarchical Scheduling, PFC, and Classifier Mapping to Interfaces\) | 263](#)

To verify that the PFC across Layer 3 interfaces configuration has been created and is operating properly, perform these tasks:

Verifying the Interface Configuration

Purpose

Verify that the Layer 2 Ethernet interfaces, Layer 3 IP interfaces, IRB interfaces, and VLAN interfaces have been created on the switch and are correctly configured.

Action

Display the switch interface configuration using the `show configuration interfaces` command:

```
user@switch> show configuration interfaces
xe-0/0/20 {
  unit 0 {
    family ethernet-switching {
      interface-mode trunk;
      vlan {
```

```

        members vlan105;
    }
}
}
xe-0/0/21 {
    unit 0 {
        family ethernet-switching {
            interface-mode trunk;
            vlan {
                members vlan106;
            }
        }
    }
}
xe-0/0/40 {
    vlan-tagging;
    unit 0 {
        vlan-id 103;
        family inet {
            address 100.103.1.2/24;
        }
    }
}
xe-0/0/41 {
    vlan-tagging;
    unit 0 {
        vlan-id 104;
        family inet {
            address 100.104.1.2/24;
        }
    }
}
irb {
    unit 105 {
        family inet {
            address 100.105.1.1/24;
        }
    }
    unit 106 {
        family inet {
            address 100.106.1.1/24;
        }
    }
}

```



```

    }
}
vlan {
    unit 105 {
        family inet {
            address 100.105.1.1/24;
        }
    }
    unit 106 {
        family inet {
            address 100.106.1.1/24;
        }
    }
}
}

```

Meaning

The `show configuration interfaces` command displays all of the interfaces configured on the switch. The command output shows that:

- Interfaces `xe-0/0/20` and `xe-0/0/21` are Ethernet interfaces (family `ethernet-switching`) in trunk interface mode. Interface `xe-0/0/20` is a member of VLAN `vlan105`, and interface `xe-0/0/21` is a member of VLAN `vlan106`.
- Interfaces `xe-0/0/40` and `xe-0/0/41` are IP interfaces (family `inet`) with VLAN tagging enabled. Interface `xe-0/0/40` has an IP address of `100.103.1.2/24` and a VLAN ID of `103`. Interface `xe-0/0/41` has an IP address of `100.104.1.2/24` and a VLAN ID of `104`.
- Two IRB interfaces are configured, IRB unit `105` with an IP address of `100.105.1.1/24` and IRB unit `106` with an IP address of `100.106.1.1/24`.
- Two VLAN interfaces are configured, VLAN unit `105` with an IP address of `100.105.1.1/24` (for IRB interface unit `105`) and VLAN unit `106` with an IP address of `100.106.1.1/24` (for IRB interface unit `106`).

Verifying the VLAN Configuration

Purpose

Verify that VLANs have been created on the switch and are correctly configured.

Action

Display the VLAN configuration using the `show configuration vlans` command:

```
user@switch> show configuration vlans
vlan105 {
    vlan-id 105;
    l3-interface irb.105;
}
vlan106 {
    vlan-id 106;
    l3-interface irb.106;
}
```

Meaning

The `show configuration vlans` command displays all of the VLANs configured on the switch. The command output shows that:

- VLAN `vlan105` has been configured with VLAN ID 105 on IRB interface `irb.105`.
- VLAN `vlan106` has been configured with VLAN ID 106 on IRB interface `irb.106`.

Verifying the PFC Configuration (Congestion Notification Profile)

Purpose

Verify that PFC has been enabled on the correct IEEE 802.1p code points (priorities) in the CNP.

Action

Display the PFC configuration using the `show configuration class-of-service congestion-notification-profile` command:

```
user@switch> show configuration class-of-service congestion-notification-profile
lossless-cnp {
    input {
        ieee-802.1 {
            code-point 011 {
                pfc;
            }
        }
    }
}
```

```

        code-point 100 {
            pfc;
        }
    }
}

```

Meaning

The `show configuration class-of-service congestion-notification-profile` command displays all of the CNPs configured on the switch. The command output shows that:

- The CNP named `lossless-cnp` is configured on the switch.
- The CNP `lossless-cnp` enables PFC on IEEE 802.1p code points 100 and 100.

Verify the Forwarding Class Configuration

Purpose

Verify that the two lossless forwarding classes and the best-effort forwarding class have been configured on the switch.

Action

Display the forwarding class configuration using the `show configuration class-of-service forwarding-classes` command:

```

user@switch> show configuration class-of-service forwarding-classes
class lossless-3 queue-num 3 no-loss;
class lossless-4 queue-num 4 no-loss;
class all-others queue-num 0;

```

Meaning

The `show configuration class-of-service forwarding-classes` command displays all of the forwarding classes configured on the switch (default forwarding classes are not displayed). The command output shows that:

- Forwarding class `lossless-3` is mapped to queue 3 and is configured as a lossless forwarding class (the `no-loss` attribute is applied)

- Forwarding class `lossless-4` is mapped to queue 4 and is configured as a lossless forwarding class (the `no-loss` attribute is applied)
- Forwarding class `all-others` is mapped to queue 0. It is not a lossless forwarding class (the `no-loss` attribute is not applied).

Verifying the Classifier Configuration

Purpose

Verify that the IEEE 802.1p classifier has been configured on the switch.

Action

Display the classifier configuration using the `show configuration class-of-service classifiers` command:

```
user@switch> show configuration class-of-service classifiers
ieee-802.1 lossless-3-4-ieee {
    forwarding-class lossless-3 {
        loss-priority low code-points 011;
    }
    forwarding-class lossless-4 {
        loss-priority low code-points 100;
    }
}
```

Meaning

The `show configuration class-of-service classifiers` command displays all of the classifiers configured on the switch. The command output shows that the Layer 2 IEEE 802.1p classifier `lossless-3-4-ieee` classifies traffic with the code point 011 into the `lossless-3` forwarding class with a loss priority of `low`, and classifies traffic with the code point 100 into the `lossless-4` forwarding class with a loss priority of `low`.

Verifying the Interface CoS Configuration (Hierarchical Scheduling, PFC, and Classifier Mapping to Interfaces)

Purpose

Verify that the interfaces have the correct hierarchical scheduling, PFC, and classifier configurations.

NOTE: The results are from the ETS hierarchical scheduling configuration, which shows the more complex configuration. Direct port scheduling results would not show the traffic control profile or forwarding class set portions of the interface configuration, but 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. Other than that, they are the same.

Action

Display the interface CoS configuration using the `show configuration class-of-service interfaces` command:

```
user@switch> show configuration class-of-service interfaces
xe-0/0/20 {
  forwarding-class-set {
    lossless_fc_set {
      output-traffic-control-profile lossless_tcp;
    }
    all-others_fc_set {
      output-traffic-control-profile all-others_tcp;
    }
  }
  congestion-notification-profile lossless-cnp;
  unit 0 {
    classifiers {
      ieee-802.1 lossless-3-4-ieee;
    }
  }
}
xe-0/0/21 {
  forwarding-class-set {
    all-others_fc_set {
      output-traffic-control-profile all-others_tcp;
    }
    lossless_fc_set {
      output-traffic-control-profile lossless_tcp;
    }
  }
  congestion-notification-profile lossless-cnp;
  unit 0 {
    classifiers {
```

```

        ieee-802.1 lossless-3-4-ieee;
    }
}
xe-0/0/40 {
    forwarding-class-set {
        lossless_fc_set {
            output-traffic-control-profile lossless_tcp;
        }
        all-others_fc_set {
            output-traffic-control-profile all-others_tcp;
        }
    }
    congestion-notification-profile lossless-cnp;
    classifiers {
        ieee-802.1 lossless-3-4-ieee;
    }
}
xe-0/0/41 {
    forwarding-class-set {
        lossless_fc_set {
            output-traffic-control-profile lossless_tcp;
        }
        all-others_fc_set {
            output-traffic-control-profile all-others_tcp;
        }
    }
    congestion-notification-profile lossless-cnp;
    classifiers {
        ieee-802.1 lossless-3-4-ieee;
    }
}

```

Meaning

The `show configuration class-of-service interfaces` command displays all of the CoS components configured on the switch interfaces. The command output shows that:

- The configuration on Layer 2 Ethernet interfaces `xe-0/0/20` and `xe-0/0/21` includes:
 - Hierarchical scheduling—The forwarding class set `lossless_fc_set` with the traffic control profile `lossless_tcp` for the lossless traffic, and the forwarding class set `all-others_fc_set` with the traffic control profile `all-others_tcp` for the best-effort traffic are applied to both interfaces.

- PFC—The `lossless-cnp` congestion notification profile is applied to both interfaces.
- Classifiers—The Layer 2 IEEE 802.1p classifier `lossless-3-4-ieee` is applied to both interfaces.
- The configuration on Layer 3 IP interfaces `xe-0/0/40` and `xe-0/0/41` includes:
 - Hierarchical scheduling—The forwarding class set `lossless_fc_set` with the traffic control profile `lossless_tcp` for the lossless traffic, and the forwarding class set `all-others_fc_set` with the traffic control profile `all-others_tcp` for the best-effort traffic are applied to both interfaces.
 - PFC—The `lossless-cnp` congestion notification profile is applied to both interfaces.
 - Classifiers—The Layer 2 IEEE 802.1p classifier `lossless-3-4-ieee` is applied to both interfaces. Traffic that would use a DSCP or a DSCP IPv6 classifier if it were configured uses the IEEE 802.1p classifier instead. Using the IEEE 802.1p classifier allows the interface to use PFC to pause traffic during periods of congestion to prevent packet loss.

RELATED DOCUMENTATION

[Understanding PFC Functionality Across Layer 3 Interfaces](#) | 236

Understanding PFC Using DSCP at Layer 3 for Untagged Traffic

IN THIS SECTION

- [Overview of DSCP-based PFC](#) | 267
- [Limitations of DSCP-based PFC](#) | 267

Protocols such as Remote Direct Memory Access (RDMA) over converged Ethernet version 2 (RoCEv2) require lossless behavior for traffic across Layer 3 connections to Layer 2 Ethernet subnetworks. Traditionally, priority-based flow control (PFC) can be used to prevent traffic loss when congestion occurs on Layer 2 or Layer 3 interfaces for VLAN-tagged traffic by selectively pausing traffic on any of eight priorities corresponding to IEEE 802.1p code points in the VLAN headers of incoming traffic on an interface. However, *untagged* traffic—traffic without VLAN tagging—cannot be examined for IEEE 802.1p code points on which to pause traffic.

Starting in Junos OS Release 17.4R1, to support lossless traffic flow at Layer 3 for untagged traffic, we support enabling PFC for Layer 3 interfaces and Layer 2 access interfaces using Distributed Services

code point (DSCP) values in the Layer 3 IP header of incoming traffic, rather than IEEE 802.1p code point values in a Layer 2 VLAN header.

Overview of DSCP-based PFC

PFC is a data center bridging technology operating at Layer 2, and DSCP information is exchanged in IP headers at Layer 3. However, you can configure DSCP-based PFC, which preserves lossless behavior across Layer 3 network connections for untagged traffic.

PFC operates by generating pause frames for traffic identified on configured code points in incoming traffic to notify the peer to pause transmission when the link is congested. With DSCP-based PFC enabled, pause frames are triggered based on a configured 6-bit DSCP value (corresponding to decimal values 0-63) in the Layer 3 IP header of incoming traffic.

However, PFC can only send pause frames with a 3-bit PFC priority—one of 8 code points corresponding to decimal values 0-7—which, for VLAN-tagged traffic, usually corresponds to the IEEE 802.1p code points in the incoming traffic VLAN headers. Untagged traffic provides no reference for IEEE 802.1p code point values, so to trigger PFC on a DSCP value, the DSCP value must be mapped explicitly in the configuration to a PFC priority to use in the PFC pause frames sent to the peer when congestion occurs for that code point. You can map traffic on a DSCP value to a PFC priority when you define the no-loss forwarding class with which you want to classify DSCP-based PFC traffic. The forwarding class must also be mapped to an output queue with no-loss behavior.

NOTE: You cannot assign the same PFC priority to more than one forwarding class because the mapped PFC priority value is used as the forwarding class ID when DSCP-based PFC is configured.

A DSCP classifier (instead of an IEEE 802.1p classifier) is also required to specify that incoming traffic with the above-configured DSCP value belongs to the no-loss forwarding class. Any DSCP values for which DSCP-based PFC is enabled on a interface must be specified in either the default DSCP classifier or in a user-defined DSCP classifier associated with the interface.

To enable DSCP-based PFC on an interface, define an input congestion notification profile with the same DSCP value (and desired buffering parameters), and associate it with the interface.

The peer device should have a matching PFC configuration for the mapped PFC priority code points.

Limitations of DSCP-based PFC

The following are limitations of DSCP-based PFC:

- You cannot configure both DSCP-based PFC and IEEE 802.1p PFC under the same congestion notification profile, or associate both a DSCP-based congestion notification profile and an IEEE 802.1p congestion notification profile with the same interface.
- DSCP-based PFC is supported on Layer 3 interfaces and Layer 2 access interfaces for untagged traffic only. PFC behavior is unpredictable if VLAN-tagged packets are received on an interface with DSCP-based PFC enabled.
- Each no-loss forwarding class can only be associated with a unique 3-bit PFC priority value from 0 through 7.

Release History Table

Release	Description
17.4R1	Starting in Junos OS Release 17.4R1, to support lossless traffic flow at Layer 3 for untagged traffic, we support enabling PFC for Layer 3 interfaces and Layer 2 access interfaces using Distributed Services code point (DSCP) values in the Layer 3 IP header of incoming traffic, rather than IEEE 802.1p code point values in a Layer 2 VLAN header.

RELATED DOCUMENTATION

<i>Configuring DSCP-based PFC for Layer 3 Untagged Traffic</i>
Understanding CoS Classifiers 96
Understanding CoS Flow Control (Ethernet PAUSE and PFC) 220
Understanding CoS Forwarding Classes 155
Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows 194

Configuring DSCP-based PFC for Layer 3 Untagged Traffic

You can configure DSCP-based PFC to support lossless behavior for untagged traffic across Layer 3 connections to Layer 2 subnetworks for protocols such as Remote Direct Memory Access (RDMA) over converged Ethernet version 2 (RoCEv2).

With DSCP-based PFC, pause frames are generated to notify the peer that the link is congested based on a configured 6-bit Distributed Services code point (DSCP) value in the Layer 3 IP header of incoming traffic, rather than a 3-bit IEEE 802.1p code point in the Layer 2 VLAN header.

Because PFC can only send pause frames corresponding to PFC priority code points, the 6-bit configured DSCP value must be mapped to a 3-bit PFC priority to use in pause frames when DSCP-based PFC is triggered. Configuring the mapping involves mapping the PFC priority value to a no-loss

forwarding class when you map the forwarding class to a queue, defining a congestion notification profile to enable PFC on traffic with the desired DSCP value, and configuring a DSCP classifier to associate the PFC priority-mapped forwarding class (along with the loss priority) with the configured DSCP value on which to trigger PFC pause frames.

The peer device should have output PFC and a corresponding flow control queue configured to match the PFC priority configuration on the device.

To configure DSCP-based PFC:

1. Map a lossless forwarding class to a PFC priority—a 3-bit value represented in decimal form (0-7)—to use in the PFC pause frames.

You must also assign an output queue to the forwarding class with the `queue-num` option. The `no-loss` option is required in this case to support lossless behavior for DSCP-based PFC, and the `pfc-priority` statement specifies the priority value mapping, as follows:

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

2. Define an input congestion notification profile to enable PFC on traffic specified by the desired 6-bit DSCP value, and optionally configure the maximum receive unit (MRU) at this time (used to determine PFC buffer headroom space reserved for the link):

```
[edit class-of-service]
user@switch# set congestion-notification-profile name input dscp code-point code-point-bits
pfc mru mru-value
```

NOTE: You cannot configure both DSCP-based PFC and IEEE 802.1p PFC under the same congestion notification profile.

3. Set up a DSCP classifier for the configured DSCP value and no-loss forwarding class mapped in the previous steps:

```
[edit class-of-service]
user@switch# set classifiers dscp classifier-name forwarding-class class-name loss-priority
level code-points code-point-bits
```

4. Assign the classifier and congestion notification profile set up in the previous steps to an interface on which you are enabling DSCP-based PFC:

```
[edit class-of-service]
user@switch# set interfaces interface-name classifiers dscp classifier-name
user @switch# set interfaces interface-name congestion-notification-profile profile-name
```

For example, with the following sample commands configuring DSCP-based PFC for interface xe-0/0/1, PFC pause frames will be generated with PFC priority 3 when incoming traffic with DSCP value 110000 becomes congested:

```
set interfaces xe-0/0/1 unit 0 family inet address 10.1.1.2/24
set class-of-service forwarding-classes class fc1 queue-num 1 no-loss
set class-of-service forwarding-classes class fc1 pfc-priority 3
set class-of-service congestion-notification-profile dpfc-cnp input dscp code-point 110000 pfc
set class-of-service classifiers dscp dpfc forwarding-class fc1 loss-priority low code-points 110000
set class-of-service interfaces xe-0/0/1 congestion-notification-profile dpfc-cnp
set class-of-service interfaces xe-0/0/1 classifiers dscp dpfc
```

RELATED DOCUMENTATION

Understanding PFC Using DSCP at Layer 3 for Untagged Traffic

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 106](#)

[Defining CoS Forwarding Classes](#)

CoS and Host Outbound Traffic

IN THIS CHAPTER

- Understanding Host Routing Engine Outbound Traffic Queues and Defaults | 271
- Changing the Host Outbound Traffic Default Queue Mapping | 274

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 60 on page 272](#).

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 60 on page 272](#) lists the default host traffic to output queue mapping.

Table 60: 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 60: 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 | 155](#)
[Understanding CoS Forwarding Classes](#)
[Changing the Host Outbound Traffic Default Queue Mapping | 274](#)
[Example: Configuring Forwarding Classes | 174](#)

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](#) | 271

2

PART

Weighted Random Early Detection (WRED) and Explicit Congestion Notification (ECN)

[WRED and Drop Profiles](#) | 276

[Explicit Congestion Notification \(ECN\)](#) | 297

WRED and Drop Profiles

IN THIS CHAPTER

- [Understanding CoS WRED Drop Profiles | 276](#)
- [Configuring CoS WRED Drop Profiles | 284](#)
- [Example: Configuring WRED Drop Profiles | 286](#)
- [Configuring CoS Drop Profile Maps | 293](#)
- [Example: Configuring Drop Profile Maps | 293](#)

Understanding CoS WRED Drop Profiles

IN THIS SECTION

- [Drop Profile Parameters | 277](#)
- [Defining Drop Profiles on Switches Except QFX10000 | 277](#)
- [Defining Drop Profiles on QFX10000 Switches | 278](#)
- [Default Drop Profile | 279](#)
- [Packet Drop Method | 280](#)
- [Packet Drop Example for Switches Except QFX10000 | 280](#)
- [Drop Profile Maps | 281](#)
- [Congestion Prevention | 281](#)
- [Configuring a WRED Drop Profile and Applying it to an Output Queue | 282](#)
- [Drop Profiles on Explicit Congestion Notification Enabled Queues | 283](#)

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.

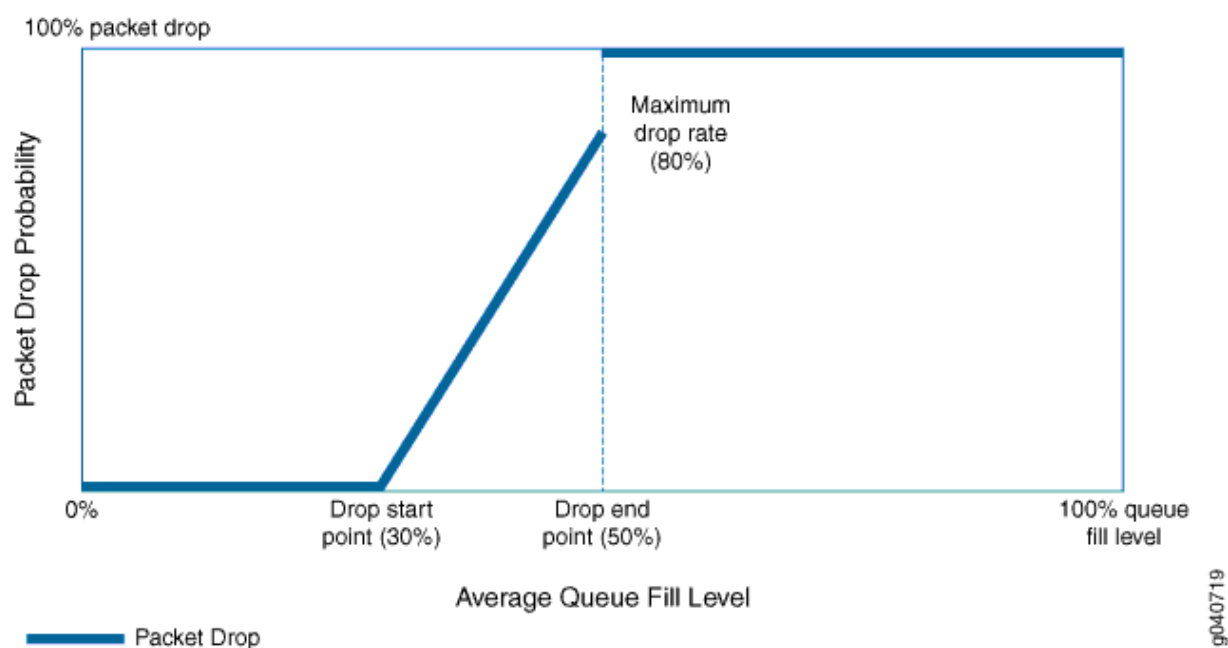
Defining Drop Profiles on Switches Except QFX10000

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 8 on page 278](#). After the drop end point, all packets drop.

Figure 8: WRED-Drop Profile Packet Drop Pattern



The thick line in [Figure 8 on page 278](#) 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 | 21](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[*Understanding CoS Explicit Congestion Notification*](#)

[Example: Configuring WRED Drop Profiles | 286](#)

[Example: Configuring Drop Profile Maps | 293](#)

[Example: Configuring Unicast Classifiers | 113](#)

[Configuring CoS WRED Drop Profiles | 284](#)

[Configuring CoS Drop Profile Maps | 293](#)

[Defining CoS BA Classifiers \(DSCP, DSCP IPv6, IEEE 802.1p\) | 106](#)

Configuring CoS WRED Drop Profiles

IN THIS SECTION

- [Drop Profiles on Switches Except QFX10000 | 285](#)
- [Drop Profiles on QFX 10000 Switches | 285](#)

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 \(ETS\) | 445](#)

[Example: Configuring WRED Drop Profiles | 286](#)

[Defining CoS Queue Schedulers | 346](#)

[Defining CoS Queue Schedulers for Port Scheduling | 382](#)

[Configuring CoS Drop Profile Maps | 293](#)

[Understanding CoS WRED Drop Profiles | 276](#)

Example: Configuring WRED Drop Profiles

IN THIS SECTION

- [Requirements | 287](#)
- [Overview | 287](#)
- [Configuring WRED Drop Profiles on Switches Except QFX10000 | 288](#)
- [Configuring WRED Drop Profiles on QFX10000 Switches | 291](#)

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 | 290](#)

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.

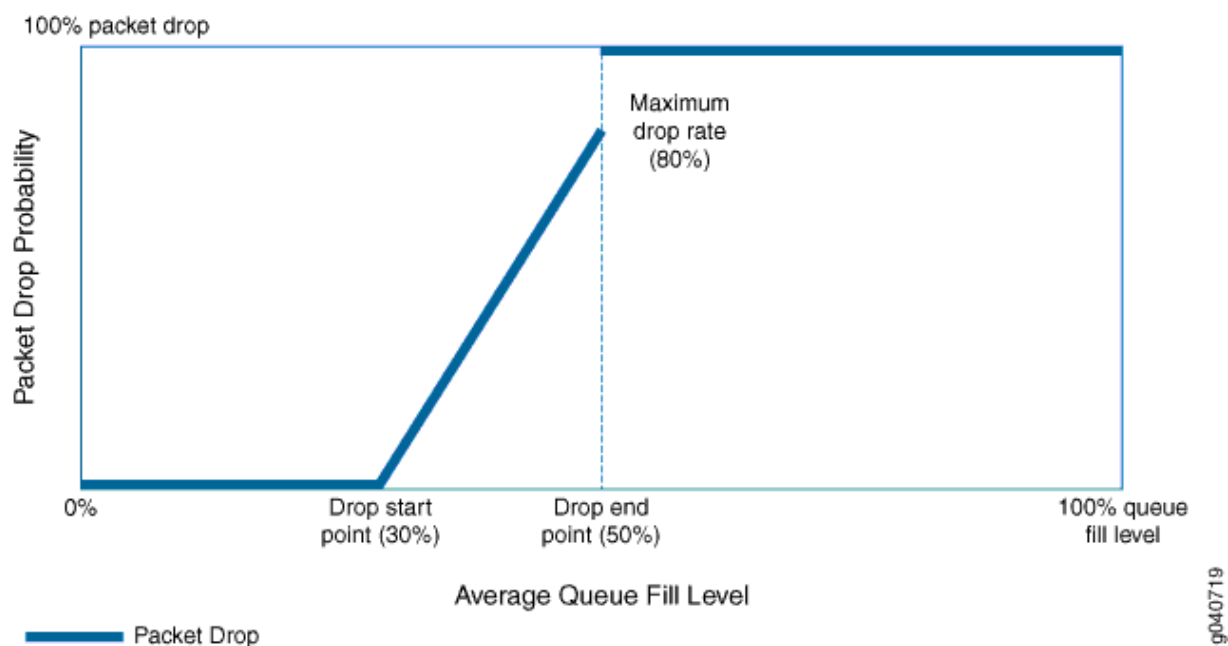
You set the drop start point and the drop end point by specifying two queue fill level percentage values. The first value is the drop start point and the second value is the drop end point.

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 9 on page 289](#) 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 9: 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 9 on page 289](#). 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 | 290](#)

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 | 292](#)

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 | 292](#)

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 | 295](#)
- [Overview | 295](#)
- [Verification | 295](#)

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 | 296](#)

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 \(ETS\) | 445](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Queue Schedulers for Port Scheduling | 386](#)

[Example: Configuring WRED Drop Profiles | 286](#)

[Configuring CoS Drop Profile Maps | 293](#)

[Understanding CoS WRED Drop Profiles | 276](#)

Explicit Congestion Notification (ECN)

IN THIS CHAPTER

- [Understanding CoS Explicit Congestion Notification | 297](#)
- [Example: Configuring ECN | 307](#)
- [Data Center Quantized Congestion Notification \(DCQCN\) | 314](#)

Understanding CoS Explicit Congestion Notification

IN THIS SECTION

- [How ECN Works | 298](#)
- [WRED Drop Profile Control of ECN Thresholds | 303](#)
- [Support, Limitations, and Notes | 306](#)

Explicit congestion notification (ECN) enables end-to-end congestion notification between two endpoints on TCP/IP based networks. The two endpoints are an ECN-enabled sender and an ECN-enabled receiver. 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.

ECN is disabled by default. Normally, you enable ECN only on queues that handle best-effort traffic because other traffic types use different methods of congestion notification—lossless traffic uses priority-based flow control (PFC) and strict-high priority traffic receives all of the port bandwidth it requires up to the point of a configured maximum rate.

You enable ECN on individual output queues (as represented by forwarding classes) by enabling ECN in the queue scheduler configuration, mapping the scheduler to forwarding classes (queues), and then applying the scheduler to interfaces.

NOTE: For ECN to work on a queue, you must also apply a weighted random early detection (WRED) packet drop profile to the queue.

How ECN Works

Without ECN, switches respond to network congestion by dropping TCP/IP packets. Dropped packets signal the network that congestion is occurring. Devices on the IP network respond to TCP packet drops by reducing the packet transmission rate to allow the congestion to clear. However, the packet drop method of congestion notification and management has some disadvantages. For example, packets are dropped and must be retransmitted. Also, bursty traffic can cause the network to reduce the transmission rate too much, resulting in inefficient bandwidth utilization.

Instead of dropping packets to signal network congestion, ECN marks packets to signal network congestion, without dropping the packets. For ECN to work, all of the switches in the path between two ECN-enabled endpoints must have ECN enabled. ECN is negotiated during the establishment of the TCP connection between the endpoints.

ECN-enabled switches determine the queue congestion state based on the WRED packet drop profile configuration applied to the queue, so each ECN-enabled queue must also have a WRED drop profile. If a queue fills to the level at which the WRED drop profile has a packet drop probability greater than zero (0), the switch might mark a packet as experiencing congestion. Whether or not a switch marks a packet as experiencing congestion is the same probability as the drop probability of the queue at that fill level.

ECN communicates whether or not congestion is experienced by marking the two least-significant bits in the differentiated services (DiffServ) field in the IP header. The most significant six bits in the DiffServ field contain the Differentiated Services Code Point (DSCP) bits. The state of the two ECN bits signals whether or not the packet is an ECN-capable packet and whether or not congestion has been experienced.

ECN-capable senders mark packets as ECN-capable. If a sender is not ECN-capable, it marks packets as not ECN-capable. If an ECN-capable packet experiences congestion at the egress queue of a switch, the switch marks the packet as experiencing congestion. When the packet reaches the ECN-capable receiver (destination endpoint), the receiver echoes the congestion indicator to the sender (source endpoint) by sending a packet marked to indicate congestion.

After receiving the congestion indicator from the receiver, the source endpoint reduces the transmission rate to relieve the congestion. This is similar to the result of TCP congestion notification and management, but instead of dropping the packet to signal network congestion, ECN marks the packet

and the receiver echoes the congestion notification to the sender. Because the packet is not dropped, the packet does not need to be retransmitted.

ECN Bits in the DiffServ Field

The two ECN bits in the DiffServ field provide four codes that determine if a packet is marked as an ECN-capable transport (ECT) packet, meaning that both endpoints of the transport protocol are ECN-capable, and if there is congestion experienced (CE), as shown in [Table 61 on page 299](#):

Table 61: ECN Bit Codes

ECN Bits (Code)	Meaning
00	Non-ECT—Packet is marked as not ECN-capable
01	ECT(1)—Endpoints of the transport protocol are ECN-capable
10	ECT(0)—Endpoints of the transport protocol are ECN-capable
11	CE—Congestion experienced

Codes 01 and 10 have the same meaning: the sending and receiving endpoints of the transport protocol are ECN-capable. There is no difference between these codes.

End-to-End ECN Behavior

After the sending and receiving endpoints negotiate ECN, the sending endpoint marks packets as ECN-capable by setting the DiffServ ECN field to ECT(1) (01) or ECT(0) (10). Every intermediate switch between the endpoints must have ECN enabled or it does not work.

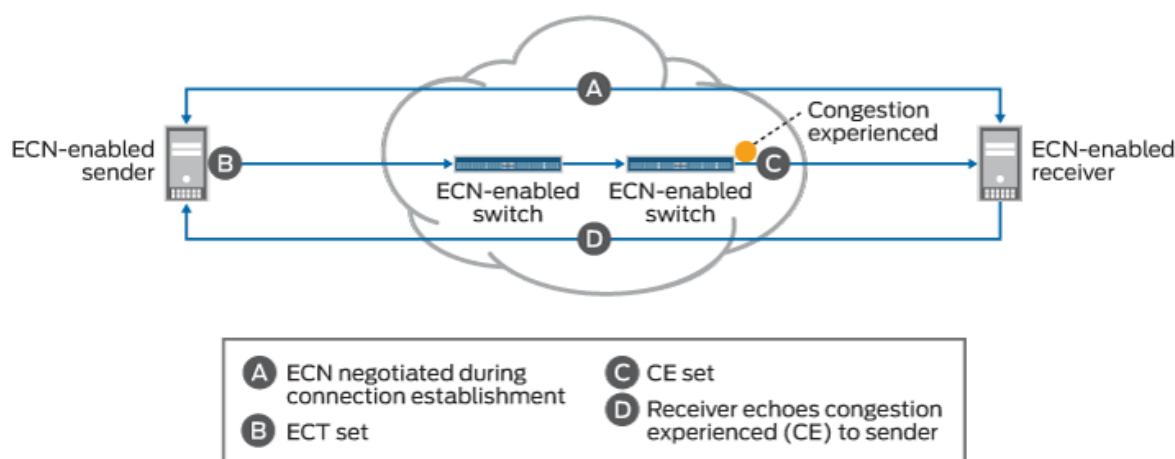
When a packet traverses a switch and experiences congestion at an output queue that uses the WRED packet drop mechanism, the switch marks the packet as experiencing congestion by setting the DiffServ ECN field to CE (11). Instead of dropping the packet (as with TCP congestion notification), the switch forwards the packet.

NOTE: At the egress queue, the WRED algorithm determines whether or not a packet is drop eligible based on the queue fill level (how full the queue is). If a packet is drop eligible and marked as ECN-capable, the packet can be marked CE and forwarded. If a packet is drop eligible and is

not marked as ECN-capable, it might be dropped. See ["WRED Drop Profile Control of ECN Thresholds"](#) on page 303 for more information about the WRED algorithm.

When the packet reaches the receiver endpoint, the CE mark tells the receiver that there is network congestion. The receiver then sends (echoes) a message to the sender that indicates there is congestion on the network. The sender acknowledges the congestion notification message and reduces its transmission rate. [Figure 10 on page 300](#) summarizes how ECN works to mitigate network congestion:

Figure 10: Explicit Congestion Notification



8042495

End-to-end ECN behavior includes:

1. The ECN-capable sender and receiver negotiate ECN capability during the establishment of their connection.
2. After successful negotiation of ECN capability, the ECN-capable sender sends IP packets with the ECT field set to the receiver.

NOTE: All of the intermediate devices in the path between the sender and the receiver must be ECN-enabled.

3. If the WRED algorithm on a switch egress queue determines that the queue is experiencing congestion and the packet is drop eligible, the switch can mark the packet as "congestion experienced" (CE) to indicate to the receiver that there is congestion on the network. If the packet has already been marked CE (congestion has already been experienced at the egress of another switch), the switch forwards the packet with CE marked.

If there is no congestion at the switch egress queue, the switch forwards the packet and does not change the ECT-enabled marking of the ECN bits, so the packet is still marked as ECN-capable but not as experiencing congestion.

On QFX5210, QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, packets that are not marked as ECN-capable (ECT, 00) are treated according to the WRED drop profile configuration and might be dropped during periods of congestion.

On QFX10000 switches, the switch uses the tail-drop algorithm to drop packets that are marked ECT (00) 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.)

4. The receiver receives a packet marked CE to indicate that congestion was experienced along the congestion path.
5. The receiver echoes (sends) a packet back to the sender with the ECE bit (bit 9) marked in the flag field of the TCP header. The ECE bit is the ECN echo flag bit, which notifies the sender that there is congestion on the network.
6. The sender reduces the data transmission rate and sends a packet to the receiver with the CWR bit (bit 8) marked in the flag field of the TCP header. The CWR bit is the congestion window reduced flag bit, which acknowledges to the receiver that the congestion experienced notification was received.
7. When the receiver receives the CWR flag, the receiver stops setting the ECE bit in replies to the sender.

[Table 62 on page 301](#) summarizes the behavior of traffic on ECN-enabled queues.

Table 62: Traffic Behavior on ECN-Enabled Queues

Incoming IP Packet Marking of ECN Bits	ECN Configuration on the Output Queue	Action if WRED Algorithm Determines Packet is Drop Eligible	Outgoing Packet Marking of ECN Bits
Non-ECT (00)	Does not matter	Drop (QFX5210, QFX5200, QFX5100, EX4600, QFX3500, QFX3600, QFabric systems). Tail drop occurs when queue reaches maximum fullness because no WRED drop probability is applied (QFX10000 switches).	No ECN bits marked

Table 62: Traffic Behavior on ECN-Enabled Queues (Continued)

Incoming IP Packet Marking of ECN Bits	ECN Configuration on the Output Queue	Action if WRED Algorithm Determines Packet is Drop Eligible	Outgoing Packet Marking of ECN Bits
ECT (10 or 01)	ECN disabled	Drop	Packet dropped—no ECN bits marked
ECT (10 or 01)	ECN enabled	Do not drop. Mark packet as experiencing congestion (CE, bits 11).	Packet marked ECT (11) to indicate congestion
CE (11)	ECN disabled	Drop	Packet dropped—no ECN bits marked
CE (11)	ECN enabled	Do not drop. Packet is already marked as experiencing congestion, forward packet without changing the ECN marking.	Packet marked ECT (11) to indicate congestion

When an output queue is not experiencing congestion as defined by the WRED drop profile mapped to the queue, all packets are forwarded, and no packets are dropped.

ECN Compared to PFC and Ethernet PAUSE

ECN is an end-to-end network congestion notification mechanism for IP traffic. Priority-based flow control (PFC) (IEEE 802.1Qbb) and Ethernet PAUSE (IEEE 802.3X) are different types of congestion management mechanisms.

ECN requires that an output queue must also have an associated WRED packet drop profile. Output queues used for traffic on which PFC is enabled should not have an associated WRED drop profile. Interfaces on which Ethernet PAUSE is enabled should not have an associated WRED drop profile.

PFC is a peer-to-peer flow control mechanism to support lossless traffic. PFC enables connected peer devices to pause flow transmission during periods of congestion. PFC enables you to pause traffic on a specified type of flow on a link instead of on all traffic on a link. For example, you can (and should) enable PFC on lossless traffic classes such as the `fcoe` forwarding class. Ethernet PAUSE is also a peer-to-peer flow control mechanism, but instead of pausing only specified traffic flows, Ethernet PAUSE pauses all traffic on a physical link.

With PFC and Ethernet PAUSE, the sending and receiving endpoints of a flow do not communicate congestion information to each other across the intermediate switches. Instead, PFC controls flows

between two PFC-enabled peer devices (for example, switches) that support data center bridging (DCB) standards. PFC works by sending a pause message to the connected peer when the flow output queue becomes congested. Ethernet PAUSE simply pauses all traffic on a link during periods of congestion and does not require DCB.

PFC works this way: if a switch output queue fills to a certain threshold, the switch sends a PFC pause message to the connected peer device that is transmitting data. The pause message tells the transmitting switch to pause transmission of the flow. When the congestion clears, the switch sends another PFC message to tell the connected peer to resume transmission. (If the output queue of the transmitting switch also reaches a certain threshold, that switch can in turn send a PFC pause message to the connected peer that is transmitting to it. In this way, PFC can propagate a transmission pause back through the network.)

See ["Understanding CoS Flow Control \(Ethernet PAUSE and PFC\)" on page 220](#) for more information. For QFX5100 and EX4600 switches only, you can also refer to ["Understanding PFC Functionality Across Layer 3 Interfaces" on page 236](#).

WRED Drop Profile Control of ECN Thresholds

You apply WRED drop profiles to forwarding classes (which are mapped to output queues) to control how the switch marks ECN-capable packets. A scheduler map associates a drop profile with a scheduler and a forwarding class, and then you apply the scheduler map to interfaces to implement the scheduling properties for the forwarding class on those interfaces.

Drop profiles define queue fill level (the percentage of queue fullness) and drop probability (the percentage probability that a packet is dropped) pairs. When a queue fills to a specified level, traffic that matches the drop profile has the drop probability paired with that fill level. When you configure a drop profile, you configure pairs of fill levels and drop probabilities to control how packets drop at different levels of queue fullness.

The first fill level and drop probability pair is the drop start point. Until the queue reaches the first fill level, packets are not dropped. When the queue reaches the first fill level, packets that exceed the fill level have a probability of being dropped that equals the drop probability paired with the fill level.

The last fill level and drop probability pair is the drop end point. When the queue reaches the last fill level, all packets are dropped unless they are configured for ECN.

NOTE: Lossless queues (forwarding class configured with the `no-loss` packet drop attribute) and strict-high priority queues do not use drop profiles. Lossless queues use PFC to control the flow of traffic. Strict-high priority queues receive all of the port bandwidth they require up to the configured maximum bandwidth limit (scheduler `transmit-rate` on QFX10000 switches, and

shaping-rate on QFX5210, QFX5200, QFX5100, QFX3500, QFX3600, and EX4600 switches, and QFabric systems).

Different switches support different amounts of fill level/drop probability pairs in drop profiles. For example, QFX10000 switches support 32 fill level/drop probability pairs, so there can be as many as 30 intermediate fill level/drop probability pairs between the drop start and drop endpoints. QFX5210, QFX5200, QFX5100, QFX3500, QFX3600, and EX4600 switches, and QFabric systems support two fill level/drop probability pairs—by definition, the two pairs you configure on these switches are the drop start and drop end points.

NOTE: Do not configure the last fill level as 100 percent.

The drop profile configuration affects ECN packets as follows:

- Drop start point—ECN-capable packets might be marked as congestion experienced (CE).
- Drop end point—ECN-capable packets are always marked CE.

As a queue fills from the drop start point to the drop end point, the probability that an ECN packet is marked CE is the same as the probability that a non-ECN packet is dropped if you apply the drop profile to best-effort traffic. As the queue fills, the probability of an ECN packet being marked CE increases, just as the probability of a non-ECN packet being dropped increases when you apply the drop profile to best-effort traffic.

At the drop end point, all ECN packets are marked CE, but the ECN packets are not dropped. When the queue fill level exceeds the drop end point, all ECN packets are marked CE. (At this point on QFX5210, QFX5200, QFX5100, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, all non-ECN packets are dropped.) ECN packets (and all other packets) are tail-dropped if the queue fills completely.

To configure a WRED packet drop profile and apply it to an output queue (using hierarchical scheduling on switches that support ETS):

1. Configure a drop profile using 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`.
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. 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. Associate the traffic control profile with an interface using the statement `set class-of-service interface 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, including the enable ECN attribute) 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.

Starting in Release 15.1, you can configure a WRED packet drop profile and apply it to an output queue on switches that support port scheduling (ETS hierarchical scheduling is either not supported or not used). To configure a WRED packet drop profile and apply it to an output queue on switches that support port scheduling (ETS hierarchical scheduling is either not supported or not used):

1. Configure a drop profile using 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) 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. 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.

Support, Limitations, and Notes

If the WRED algorithm that is mapped to a queue does not find a packet drop eligible, then the ECN configuration and ECN bits marking does not matter. The packet transport behavior is the same as when ECN is not enabled.

ECN is disabled by default. Normally, you enable ECN only on queues that handle best-effort traffic, and you do not enable ECN on queues that handle lossless traffic or strict-high priority traffic.

ECN supports the following:

- IPv4 and IPv6 packets
- Untagged, single-tagged, and double-tagged packets
- The outer IP header of IP tunneled packets (but not the inner IP header)

ECN does not support the following:

- IP packets with MPLS encapsulation
- The inner IP header of IP tunneled packets (however, ECN works on the outer IP header)
- Multicast, broadcast, and destination lookup fail (DLF) traffic
- Non-IP traffic

NOTE: On QFX10000 switches, when you enable a queue for ECN and apply a WRED drop profile to the queue, the WRED drop profile only sets the thresholds for marking ECN traffic as experiencing congestion (CE, 11). On ECN-enabled queues, the WRED drop profile does not set drop thresholds for non-ECT (00) traffic (traffic that is not ECN-capable). Instead, the switch uses the tail-drop algorithm on traffic that is marked non-ECT on ECN-enabled queues during periods of congestion.

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.

Release History Table

Release	Description
15.1	Starting in Release 15.1, you can configure a WRED packet drop profile and apply it to an output queue on switches that support port scheduling (ETS hierarchical scheduling is either not supported or not used).

RELATED DOCUMENTATION

| *Example: Configuring ECN*

Example: Configuring ECN

IN THIS SECTION

- [Requirements | 307](#)
- [Overview | 307](#)
- [Configuration | 310](#)
- [Verification | 313](#)

This example shows how to enable explicit congestion notification (ECN) on an output queue.

Requirements

This example uses the following hardware and software components:

- One switch.
- Junos OS Release 13.2X51-D25 or later for the QFX Series or Junos OS Release 14.1X53-D20 for the OCX Series

Overview

ECN enables end-to-end congestion notification between two endpoints on TCP/IP based networks. The two endpoints are an ECN-enabled sender and an ECN-enabled receiver. 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.

A weighted random early detection (WRED) packet drop profile must be applied to the output queues on which ECN is enabled. ECN uses the WRED drop profile thresholds to mark packets when the output queue experiences congestion.

ECN reduces packet loss by forwarding ECN-capable packets during periods of network congestion instead of dropping those packets. (TCP notifies the network about congestion by dropping packets.) During periods of congestion, ECN marks ECN-capable packets that egress from congested queues. When the receiver receives an ECN packet that is marked as experiencing congestion, the receiver

echoes the congestion state back to the sender. The sender then reduces its transmission rate to clear the congestion.

ECN is disabled by default. You can enable ECN on best-effort traffic. ECN should not be enabled on lossless traffic queues, which uses priority-based flow control (PFC) for congestion notification, and ECN should not be enabled on strict-high priority traffic queues.

To enable ECN on an output queue, you not only need to enable ECN in the queue scheduler, you also need to:

- Configure a WRED packet drop profile.
- Configure a queue scheduler that includes the WRED drop profile and enables ECN. (This example shows only ECN and drop profile configuration; you can also configure bandwidth, priority, and buffer settings in a scheduler.)
- Map the queue scheduler to a forwarding class (output queue) in a scheduler map.
- Starting in Junos OS 15.1, enhanced transmission selection (ETS) hierarchical scheduling is supported. If you are using enhanced transmission selection (ETS) hierarchical scheduling, add the forwarding class to a forwarding class set (priority group).
- If you are using ETS, associate the queue scheduler map with a traffic control profile (priority group scheduler for hierarchical scheduling).
- If you are using ETS, apply the traffic control profile and the forwarding class set to an interface. On that interface, the output queue uses the scheduler mapped to the forwarding class, as specified by the scheduler map attached to the traffic control profile. This enables ECN on the queue and applies the WRED drop profile to the queue.

If you are using port scheduling, apply the scheduler map to an interface. On that interface, the output queue uses the scheduler mapped to the forwarding class in the scheduler map, which enables ECN on the queue and applies the WRED drop profile to the queue.

[Table 63 on page 308](#) shows the configuration components for this example.

Table 63: Components of the ECN Configuration Example

Component	Settings
Hardware	QFX Series switch

Table 63: Components of the ECN Configuration Example (*Continued*)

Component	Settings
Drop profile (with two fill level/ drop probability pairs)	Name: be-dp Drop start fill level: 30 percent Drop end fill level: 75 percent Drop probability at drop start (minimum drop rate): 0 percent Drop probability at drop end (maximum drop rate): 80 percent
Scheduler	Name: be-sched ECN: enabled Drop profile: be-dp Transmit rate: 25% Buffer size: 25% Priority: low
Scheduler map	Name: be-map Forwarding class: best-effort Scheduler: be-sched NOTE: By default, the best-effort forwarding class is mapped to output queue 0.
Forwarding class set (ETS only)	Name: be-pg Forwarding class: best-effort (queue 0)
Traffic control profile (ETS only)	Name: be-tcp Scheduler map: be-map
Interface (ETS only)	Name: xe-0/0/20 Forwarding class set: be-pg (Output) traffic control profile: be-tcp
Interface (port scheduling only)	Name: xe-0/0/20

NOTE: Only switches that support ETS hierarchical scheduling support forwarding class set and traffic control profile configuration. Direct port scheduling does not use the hierarchical scheduling structure.

NOTE: On QFX5100, EX4600, QFX3500, and QFX3600 switches, and on QFabric systems, the WRED drop profile also controls packet drop behavior for traffic that is not ECN-capable (packets marked non-ECT, ECN bit code 00).

On QFX10000 switches, when ECN is enabled on a queue, the WRED drop profile only sets the ECN thresholds, it does not control packet drop on non-ECN packets. On ECN-enabled queues, QFX10000 switches use the tail-drop algorithm on non-ECN packets during periods of congestion. If you do not enable ECN, then the queue uses the WRED packet drop mechanism.

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 310](#)
- [Configuring ECN | 311](#)

CLI Quick Configuration

To quickly configure the drop profile, scheduler with ECN enabled, and to map the scheduler to an output queue on an 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.

ETS Quick Configuration

```
[edit class-of-service]
set drop-profile be-dp interpolate fill-level 30 fill-level 75 drop-probability 0 drop-
probability 80
set schedulers be-sched explicit-congestion-notification
set schedulers be-sched drop-profile-map loss-priority low protocol any drop-profile be-dp
set schedulers be-sched transmit-rate percent 25
set schedulers be-sched buffer-size percent 25
```

```

set schedulers be-sched priority low
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set forwarding-class-sets be-pg class best-effort
set traffic-control-profiles be-tcp scheduler-map be-map
set interfaces xe-0/0/20 forwarding-class-set be-pg output-traffic-control-profile be-tcp

```

Port Scheduling Quick Configuration (QFX10000 Switches)

```

[edit class-of-service]
set drop-profile be-dp interpolate fill-level 30 fill-level 75 drop-probability 0 drop-
probability 80
set schedulers be-sched explicit-congestion-notification
set schedulers be-sched drop-profile-map loss-priority low protocol any drop-profile be-dp
set schedulers be-sched transmit-rate percent 25
set schedulers be-sched buffer-size percent 25
set schedulers be-sched priority low
set scheduler-maps be-map forwarding-class best-effort scheduler be-sched
set interfaces xe-0/0/20 scheduler-map be-map

```

Configuring ECN

Step-by-Step Procedure

To configure ECN:

1. Configure the WRED packet drop profile be-dp. This example uses a drop start point of 30 percent, a drop end point of 75 percent, a minimum drop rate of 0 percent, and a maximum drop rate of 80 percent:

```

[edit class-of-service]
user@switch# set drop-profile be-dp interpolate fill-level 30 fill-level 75 drop-probability
0 drop-probability 80

```

2. Create the scheduler be-sched with ECN enabled and associate the drop profile be-dp with the scheduler:

```

[edit class-of-service]
user@switch# set schedulers be-sched explicit-congestion-notification
user@switch# set schedulers be-sched drop-profile-map loss-priority low protocol any drop-

```

```

profile be-dp
user@switch# set be-sched transmit-rate percent 25
user be-sched transmit-rate percent 25
user@switch# set be-sched buffer-size percent 25
user@switch# set be-sched buffer-size percent 25
user@switch# set be-sched priority low

```

3. Map the scheduler be-sched to the best-effort forwarding class (output queue 0) using scheduler map be-map:

```

[edit class-of-service]
user@switch# set scheduler-maps be-map forwarding-class best-effort scheduler be-sched

```

4. If you are using ETS, add the forwarding class best-effort to the forwarding class set be-pg; if you are using direct port scheduling, skip this step:

```

[edit class-of-service]
user@switch# set forwarding-class-sets be-pg class best-effort

```

5. If you are using ETS, associate the scheduler map be-map with the traffic control profile be-tcp; if you are using direct port scheduling, skip this step:

```

[edit class-of-service]
user@switch# set traffic-control-profiles be-tcp scheduler-map be-map

```

6. If you are using ETS, associate the traffic control profile be-tcp and the forwarding class set be-pg with the interface on which you want to enable ECN on the best-effort queue:

```

[edit class-of-service]
user@switch# set interfaces xe-0/0/20 forwarding-class-set be-pg output-traffic-control-
profile be-tcp

```

If you are using direct port scheduling, associate the scheduler map `be-map` with the interface on which you want to enable ECN on the best-effort queue:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/20 scheduler-map be-map
```

Verification

IN THIS SECTION

- [Verifying That ECN Is Enabled | 313](#)

Verifying That ECN Is Enabled

Purpose

Verify that ECN is enabled in the scheduler `be-sched` by showing the configuration for the scheduler map `be-map`.

Action

Display the scheduler map configuration using the operational mode command `show class-of-service scheduler-map be-map`:

```
user@switch> show class-of-service scheduler-map be-map
Scheduler map: be-map, Index: 12240

Scheduler:be-sched, Forwarding class: best-effort, Index: 115
  Transmit rate: 25 percent, Rate Limit: none, Buffer size: 25 percent,
  Buffer Limit: none, Priority: low
  Excess Priority: unspecified, Explicit Congestion Notification: enable
  Drop profiles:
    Loss priority  Protocol  Index  Name
    Low           any      3312   be-dp
    Medium-high   any      1      <default-drop-profile>
    High          any      1      <default-drop-profile>
```

Meaning

The `show class-of-service scheduler-map` operational command shows the configuration of the scheduler associated with the scheduler map and the forwarding class mapped to that scheduler. The output shows that:

- The scheduler associated with the scheduler map is `be-sched`.
- The scheduler map applies to the forwarding class `best-effort` (output queue 0).
- The scheduler `be-sched` has a transmit rate of 25 percent, a queue buffer size of 25 percent, and a drop priority of `low`.
- Explicit congestion notification state is `enable`.
- The WRED drop profile used for low drop priority traffic is `be-dp`.

Release History Table

Release	Description
15.1	Starting in Junos OS 15.1, enhanced transmission selection (ETS) hierarchical scheduling is supported.

RELATED DOCUMENTATION

| *Understanding CoS Explicit Congestion Notification*

Data Center Quantized Congestion Notification (DCQCN)

IN THIS SECTION

- [Understanding Data Center Quantized Congestion Notification \(DCQCN\) | 315](#)
- [Configuring Data Center Quantized Congestion Notification \(DCQCN\) | 316](#)

Remote Direct Memory Access (RDMA) provides the high throughput and ultra-low latency, with low CPU overhead, necessary for modern datacenter applications. RDMA is deployed using the RoCEv2 protocol, which relies on Priority-based Flow Control (PFC) to enable a drop-free network. Data Center Quantized Congestion Notification (DCQCN) is an end-to-end congestion control scheme for RoCEv2.

Starting in Junos OS Release 18.1R1, Junos OS supports DCQCN by combining Explicit Congestion Notification (ECN) and PFC to overcome the limitations of PFC to support end-to-end lossless Ethernet.

Understanding Data Center Quantized Congestion Notification (DCQCN)

Priority-based Flow Control (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. However, 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.

Explicit congestion notification (ECN) enables end-to-end congestion notification between two endpoints on TCP/IP based networks. The two endpoints are an ECN-enabled sender and an ECN-enabled receiver. 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.

Data Center Quantized Congestion Notification (DCQCN) is a combination of ECN and PFC to support end-to-end lossless Ethernet. ECN helps overcome the limitations of PFC to achieve lossless Ethernet. The idea behind DCQCN is to allow ECN to do flow control by decreasing the transmission rate when congestion starts, thereby minimizing the time PFC is triggered, which stops the flow altogether.

The correct operation of DCQCN requires balancing two conflicting requirements:

1. Ensuring PFC is not triggered too early, that is, before giving ECN a chance to send congestion feedback to slow the flow.
2. Ensuring PFC is not triggered too late, thereby causing packet loss due to buffer overflow.

There are three important parameters that need to be calculated and configured properly to achieve the above key requirements:

1. **Headroom Buffers**—A PAUSE message sent to an upstream device takes some time to arrive and take effect. To avoid packet drops, the PAUSE sender must reserve enough buffer to process any packets it may receive during this time. This includes packets that were in flight when the PAUSE was

sent, and the packets sent by the upstream device while it is processing the PAUSE message. In QFX5000 Series switches, headroom buffers are allocated on a per port per priority basis. Headroom buffers are carved out of the global shared buffer. You can control the amount of headroom buffers allocated for each port and priority using the MRU and cable length parameters in the congestion notification profile. If you see minor ingress drops even after PFC is triggered, you can eliminate those drops by increasing the headroom buffers for that port and priority combination.

2. **PFC Threshold**—This is an ingress threshold. This is the maximum size an ingress priority group can grow to before a PAUSE message is sent to the upstream device. Each PFC priority gets its own priority group at each ingress port. PFC thresholds are set per priority group at each ingress port. On QFX Series devices, there are two components in the PFC threshold—the PG MIN threshold and the PG shared threshold. Once PG MIN and PG shared thresholds are reached for a priority group, PFC is generated for that corresponding priority. The switch sends a RESUME message when the queue falls below the PFC thresholds.
3. **ECN Threshold**—This is an egress threshold. The ECN threshold is equal to the WRED start-fill-level value. Once an egress queue exceeds this threshold, the switch starts ECN marking for packets on that queue. For DCQCN to be effective, this threshold must be lower than the ingress PFC threshold to ensure PFC is not triggered before the switch has a chance to mark packets with ECN. Setting a very low WRED fill level increases ECN marking probability. For example with default shared buffer setting, a WRED start-fill-level of 10 percent ensures lossless packets are ECN marked. But with a higher fill level, the probability of ECN marking is reduced. For example, with two ingress port with lossless traffic to the same egress port and a WRED start-fill-level of 50 percent, no ECN marking will occur, because ingress PFC thresholds will be met first.

Configuring Data Center Quantized Congestion Notification (DCQCN)

To enable DCQCN, configure both ECN and PFC for a traffic flow. As an example, consider a QFX5000 Series switch between a reaction point (RP) and a notification point (NP), with et-0/0/3 as the ingress port and et-0/0/4 as the egress port.

1. Configure ECN on the egress port for a lossless flow. For example:

```
[edit class-of-service]
user@host# set drop-profiles dp1 interpolate fill-level 10 drop-probability 0 fill-level 80
drop-probability 100
user@host# set schedulers s1 drop-profile-map loss-priority any protocol any drop-profile dp1
user@host# set schedulers s1 explicit-congestion-notification
user@host# set scheduler-maps sm1 forwarding-class fcoe scheduler s1
user@host# set interfaces et-0/0/4 scheduler-map sm1
```

2. Configure PFC on the ingress port for the same lossless flow. For example:

```
[edit class-of-service]
user@host# set congestion-notification-profile cnp1 input ieee-802.1 code-point 011 pfc
user@host# set interfaces et-0/0/3 congestion-notification-profile cnp1
```

3. Configure the shared buffers. For example:

```
[edit class-of-service]
user@host# set shared-buffer ingress buffer-partition lossless percent 15
user@host# set shared-buffer ingress buffer-partition lossy percent 5
user@host# set shared-buffer ingress buffer-partition lossless-headroom percent 80
user@host# set shared-buffer egress buffer-partition lossless percent 60
user@host# set shared-buffer egress buffer-partition lossy percent 20
user@host# set shared-buffer egress buffer-partition multicast percent 20
```

4. Verify your configuration.

```
[edit class-of-service]
user@host# show
```

For example:

```
[edit class-of-service]
user@host# show
drop-profiles {
  dp1 {
    interpolate {
      fill-level [ 10 80 ];
      drop-probability [ 0 100 ];
    }
  }
}
shared-buffer {
  ingress {
    buffer-partition lossless {
      percent 15;
    }
    buffer-partition lossy {
      percent 5;
```

```

    }
    buffer-partition lossless-headroom {
        percent 80;
    }
}
egress {
    buffer-partition lossless {
        percent 60;
    }
    buffer-partition lossy {
        percent 20;
    }
    buffer-partition multicast {
        percent 20;
    }
}
}
congestion-notification-profile {
    cnp1 {
        input {
            ieee-802.1 {
                code-point 011 {
                    pfc;
                }
            }
        }
    }
}
}
interfaces {
    et-0/0/3 {
        congestion-notification-profile cnp1;
    }
    et-0/0/4 {
        scheduler-map sm1;
    }
}
scheduler-maps {
    sm1 {
        forwarding-class fcoe scheduler s1;
    }
}
}
schedulers {
    s1 {

```

```
        drop-profile-map loss-priority any protocol any drop-profile dp1;  
        explicit-congestion-notification;  
    }  
}
```

5. Save your configuration.

```
[edit class-of-service]  
user@host# commit
```

RELATED DOCUMENTATION

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

[Understanding CoS Explicit Congestion Notification | 297](#)

[Example: Configuring ECN | 307](#)

3

PART

CoS Queue Schedulers, Traffic Control Profiles, and Hierarchical Port Scheduling (ETS)

[Queue Schedulers and Scheduling Priority](#) | 321

[Port Scheduling and Shaping](#) | 368

[Troubleshooting Egress Bandwidth Issues](#) | 396

[Traffic Control Profiles and Priority Group Scheduling](#) | 401

[Hierarchical Port Scheduling \(ETS\)](#) | 438

CHAPTER 11

Queue Schedulers and Scheduling Priority

IN THIS CHAPTER

- [Understanding Default CoS Scheduling and Classification | 321](#)
- [Understanding CoS Scheduling Behavior and Configuration Considerations | 332](#)
- [Understanding CoS Output Queue Schedulers | 338](#)
- [Defining CoS Queue Schedulers | 346](#)
- [Example: Configuring Queue Schedulers | 350](#)
- [Defining CoS Queue Scheduling Priority | 358](#)
- [Example: Configuring Queue Scheduling Priority | 360](#)
- [Monitoring CoS Scheduler Maps | 365](#)

Understanding Default CoS Scheduling and Classification

IN THIS SECTION

- [Default Classification | 322](#)
- [Default Scheduling | 327](#)
- [Default DCBX Advertisement | 331](#)
- [Default Scheduling and Classification Summary | 331](#)

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 64 on page 323 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 64: 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 65 on page 324 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 65 on page 324 applies only to unicast traffic except on QFX10000 switches and NFX Series devices.

Table 65: 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
110	best-effort	low
111	best-effort	low

Table 66 on page 324 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 66 on page 324 does not apply to QFX10000 switches or NFX Series devices.

Table 66: Default IEEE 802.1 Multidestination Classifiers

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

Table 66: Default IEEE 802.1 Multidestination Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
ef1 (011)	mcast	low
af11 (100)	mcast	low
af12 (101)	mcast	low
nc1 (110)	mcast	low
nc2 (111)	mcast	low

Table 67 on page 325 shows the default mapping of DSCP code-point values to forwarding classes and loss priorities for DSCP IP and DCSP IPv6.

NOTE: Table 67 on page 325 applies only to unicast traffic except on QFX10000 switches and NFX Series devices.

Table 67: 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

Table 67: Default DSCP IP and IPv6 Classifiers (Continued)

Code Point	Forwarding Class	Loss Priority
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
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 68 on page 327](#) 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 68: 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 69 on page 328](#):

Table 69: 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 69 on page 328](#) 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 multidestination (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 (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, 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 | 26](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

[Understanding Default CoS Settings | 30](#)

[Understanding CoS Virtual Output Queues \(VOQs\) | 406](#)

[Understanding Applying CoS Classifiers and Rewrite Rules to Interfaces | 130](#)

[Understanding DCB Features and Requirements | 482](#)

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

[Example: Configuring Unicast Classifiers | 113](#)

[Example: Configuring Queue Schedulers | 350](#)

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 multidestination 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 multidestination traffic flow and a lossless unicast traffic flow (such as the default lossless `fcoe` or `no-loss` forwarding classes). Any code point used for multidestination traffic on a port should not be used to classify unicast traffic into a lossless forwarding class on the same port.

If a multidestination traffic flow and a lossless unicast traffic flow use the same code point on a port, the multidestination 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 multidestination traffic of the same code point is also paused. During periods of congestion, treating multidestination traffic the same as lossless unicast traffic can create ingress port congestion for the multidestination traffic and affect the multidestination traffic on all of the egress ports the multidestination traffic uses.

For example, the following configuration can cause ingress port congestion for the multidestination flow:

1. 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
```

2. For multidestination 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
```

3. 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 multidestination 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 | 339](#)
- [Default Schedulers | 341](#)
- [Transmit Rate \(Minimum Guaranteed Bandwidth\) | 342](#)
- [Sharing Extra Bandwidth | 343](#)
- [Shaping Rate \(Maximum Bandwidth\) | 343](#)
- [Scheduling Priority | 344](#)
- [Scheduler Drop-Profile Maps | 344](#)
- [Buffer Size | 344](#)
- [Explicit Congestion Notification | 345](#)
- [Scheduler Maps | 345](#)

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 70 on page 339](#) provides a quick reference to the scheduler components you can configure to determine the bandwidth properties of output queues (forwarding classes), and [Table 71 on page 341](#) provides a quick reference to some related scheduling configuration components.

Table 70: 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" on page 684.</p>

Table 70: Output Queue Scheduler Components *(Continued)*

Output Queue Scheduler Component	Description
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" on page 293.</p>
Explicit congestion notification	<p>Enables explicit congestion notification (ECN) on the queue.</p> <p>See <i>Understanding CoS Explicit Congestion Notification</i>.</p>
Priority	<p>Sets the scheduling priority applied to the queue.</p> <p>See "Defining CoS Queue Scheduling Priority" on page 358.</p>
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)" on page 428.</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" on page 417.</p>

Table 71: Other Scheduling Components

Other Scheduling Components	Description
Forwarding class	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.)
Output queue	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.
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" on page 327](#).

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" on page 417](#).

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 \(Maximum Bandwidth\)" on page 428](#).

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" on page 358](#).

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" on page 684](#).

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 | 21](#)

[Understanding CoS Priority Group Scheduling | 403](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

[Understanding CoS Buffer Configuration | 684](#)

[*Understanding CoS Explicit Congestion Notification*](#)

[Understanding CoS Scheduling Behavior and Configuration Considerations | 332](#)

[*Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports*](#)

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

[Configuring CoS Drop Profile Maps | 293](#)

[Defining CoS Queue Scheduling Priority | 358](#)

[Understanding CoS Priority Group Shaping and Queue Shaping \(Maximum Bandwidth\) | 428](#)

[Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth | 417](#)

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 \(ETS\)" on page 445](#) 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 \(ETS\)" on page 445](#) 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 \(ETS\) | 445](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)

[Example: Configuring Maximum Output Bandwidth | 431](#)

Example: Configuring ECN

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Defining CoS Queue Scheduling Priority | 358](#)

[Configuring CoS WRED Drop Profiles | 284](#)

[Monitoring CoS Scheduler Maps | 365](#)

[Understanding CoS Output Queue Schedulers | 338](#)

[Understanding CoS Priority Group Scheduling | 403](#)

[Understanding CoS Buffer Configuration | 684](#)

Understanding CoS Explicit Congestion Notification

Example: Configuring Queue Schedulers

IN THIS SECTION

- [Requirements | 352](#)
- [Overview | 352](#)
- [Verification | 356](#)

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 \(ETS\)](#)" on page 445 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 72 on page 355](#) shows the configuration components for this example.

Table 72: 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 | 356](#)
- [Verifying the Scheduler Map Configuration | 356](#)
- [Verifying That the Scheduler Is Associated with the Interface | 357](#)

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 \(ETS\) | 445](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)

[Example: Configuring Maximum Output Bandwidth | 431](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring WRED Drop Profiles | 286](#)

Example: Configuring ECN

[Defining CoS Queue Schedulers | 346](#)

[Monitoring CoS Scheduler Maps | 365](#)

[Understanding CoS Output Queue Schedulers | 338](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

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 CoS Buffer Configuration | 684](#)

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 | 360](#)

[Monitoring CoS Scheduler Maps | 365](#)

Example: Configuring Queue Scheduling Priority

IN THIS SECTION

- [Requirements | 361](#)
- [Overview | 361](#)
- [Verification | 363](#)

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 73 on page 362 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 `nl-sched`. One scheduler map, `schedmap1`, associates the schedulers to the queues.

Table 73: Components of the Queue Scheduler Priority Configuration Example

Component	Settings
Hardware	One switch

Table 73: Components of the Queue Scheduler Priority Configuration Example (Continued)

Component	Settings
Schedulers	fcoe-sched for FCoE traffic nl-sched for no-loss traffic
Priority	low for FCoE traffic low for no-loss traffic
Scheduler map	schedmap1: FCoE mapping: scheduler fcoe-sched to forwarding class fcoe No-loss mapping: scheduler nl-sched to forwarding class no-loss

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 | 364](#)
- [Verifying the Scheduler-to-Forwarding-Class Mapping | 364](#)

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	358
Monitoring CoS Scheduler Maps	365

Monitoring CoS Scheduler Maps

IN THIS SECTION

- Purpose | 365
- Action | 365
- Meaning | 365

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 74 on page 366](#) summarizes key output fields for CoS scheduler maps.

Table 74: Summary of Key CoS Scheduler Maps Output Fields

Field	Values
Scheduler map	Name of a scheduler map that maps forwarding classes to schedulers.
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"> 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. low—Packets in this queue are transmitted after packets in the strict-high queue.
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.

Table 74: Summary of Key CoS Scheduler Maps Output Fields *(Continued)*

Field	Values
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.
Name	Name of the drop profile.

Port Scheduling and Shaping

IN THIS CHAPTER

- [Understanding CoS Port Schedulers | 368](#)
- [Defining CoS Queue Schedulers for Port Scheduling | 382](#)
- [Example: Configuring Queue Schedulers for Port Scheduling | 386](#)
- [CoS Port Shaping | 393](#)

Understanding CoS Port Schedulers

IN THIS SECTION

- [Queue Scheduling Components | 369](#)
- [Default Schedulers | 371](#)
- [Scheduling Priority | 373](#)
- [Bandwidth Scheduling | 375](#)
- [Scheduler Drop-Profile Maps | 379](#)
- [Buffer Size | 380](#)
- [Explicit Congestion Notification | 381](#)
- [Scheduler Maps | 382](#)

Port scheduling defines the class-of-service (CoS) properties of output queues. You configure CoS properties in a scheduler, then map the scheduler to a forwarding class. Forwarding classes are in turn mapped to output queues. Classifiers map incoming traffic into forwarding classes based on IEEE 802.1p, DSCP, or EXP 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 scheduling priority of the queue, and the weighted

random early detection (WRED) drop profiles associated with the queue to control packet drop during periods of congestion.

Scheduler maps map schedulers to forwarding classes. The output queue mapped to a forwarding class receives the port resources and properties defined in the scheduler mapped to that forwarding class. You apply a scheduler map to an interface to apply queue scheduling to a port. You can associate different scheduler maps with different interfaces to configure port-specific scheduling for forwarding classes (output queues).

NOTE: Port scheduling is simpler to configure than enhanced transmission selection (ETS) two-tier hierarchical port scheduling. Port scheduling allocates port bandwidth to output queues directly, instead of allocating port bandwidth to output queues through a scheduling hierarchy. While port scheduling is simpler, ETS is more flexible. ETS allocates port bandwidth in a two-tier hierarchy:

- Port bandwidth is first allocated to a priority group using the CoS properties defined in a traffic control profile. A priority group is a group of forwarding classes (which are mapped to output queues) that require similar CoS treatment.
- Priority group bandwidth is allocated to the output queues (which are mapped to forwarding classes) using the properties defined in the output queue scheduler.

NOTE: When you configure bandwidth for a queue, 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, consider the preamble and the IFG as well as the data in the calculations.

Queue Scheduling Components

[Table 75 on page 369](#) provides a quick reference to the scheduler components you can configure to determine the bandwidth properties of output queues (forwarding classes).

Table 75: Output Queue Scheduler Components

Output Queue Scheduler Component	Description
Buffer size	Sets the size of the queue buffer.

Table 75: Output Queue Scheduler Components (Continued)

Output Queue Scheduler Component	Description
Drop profile map	<p>Maps a drop profile to a packet 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 packet loss priority to which a drop profile applies.
Excess rate	<p>Sets the percentage of extra bandwidth (bandwidth that is not used by other queues) a queue can receive. If not set, the switch uses the transmit rate to determine how much extra bandwidth the queue can use. Extra bandwidth is the bandwidth remaining after all guaranteed bandwidth requirements are met.</p>
Explicit congestion notification	<p>Enables explicit congestion notification (ECN) on the queue.</p>
Priority	<p>Sets the scheduling priority applied to the queue.</p>
Transmit rate	<p>Sets the minimum guaranteed bandwidth on low and high priority queues. By default, if you do not configure an excess rate, extra bandwidth is shared among queues in proportion to the transmit rate of each queue.</p> <p>On strict-high priority queues, sets the amount of bandwidth that receives strict-high priority forwarding treatment. 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.</p> <p>If you configure two or more strict-high priority queues on a port, you must configure a transmit rate on those queues. However, we strongly recommend that you always configure a transmit rate on strict-high priority queues to prevent them from starving other queues.</p>

[Table 76 on page 371](#) provides a quick reference to some related scheduling configuration components.

Table 76: Related Scheduling Components

Related Scheduling Components	Description
Forwarding class	Maps traffic classified into the forwarding class at the switch ingress 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 map to output queues. This mapping determines the output queue each class of traffic uses on the switch egress interfaces.)
Output queue (virtual output queue)	Output queues are virtual, and are comprised of the physical buffers on the ingress pipeline of each Packet Forwarding Engine (PFE) chip to store traffic for every egress port. Every output queue on an egress port has buffer storage space on every ingress pipeline on all of the PFE chips on the switch. The mapping of ingress pipeline storage space to output queues is 1-to-1, so each output queue receives buffer space on each ingress pipeline. See "Understanding CoS Virtual Output Queues (VOQs) on QFX10000 Switches" on page 406 for more information.
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).

Default Schedulers

If you do not configure CoS, the switch uses its default settings. Each forwarding class requires a scheduler to set the CoS properties of the forwarding class and its output queue. The default configuration has four forwarding classes: best-effort (queue 0), fcoe (queue 3), no-loss (queue 4), and network-control (queue 7). Each default forwarding class is mapped to a default scheduler. You can use the default schedulers or you can define new schedulers for these four forwarding classes. For explicitly configured forwarding classes, you must explicitly configure a queue scheduler to allocate CoS resources to the traffic mapped to each forwarding class.

[Table 77 on page 372](#) shows the default queue schedulers.

Table 77: Default Scheduler Configuration

Default Scheduler and Queue Number	Transmit Rate (Guaranteed Minimum Bandwidth)	Rate Shaping (Maximum Bandwidth)	Excess Bandwidth Sharing	Priority	Buffer Size
best-effort forwarding class scheduler (queue 0)	15%	None	15%	low	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)	15%	None	15%	low	15%

NOTE: By default, the minimum guaranteed bandwidth (transmit rate) determines the amount of excess (extra) bandwidth a queue can share. Extra bandwidth is allocated to queues in proportion to the transmit rate of each queue. You can configure bandwidth sharing (excess rate) to override the default setting and configure the excess bandwidth percentage independently of the transmit rate.

By default, only the four default schedulers shown in [Table 77 on page 372](#) 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 default bandwidth of a default queue.) If a forwarding class does not transport traffic, the bandwidth allocated to that forwarding class is available to other forwarding classes. Unicast and multidestination (multicast, broadcast, and destination lookup fail) traffic use the same forwarding classes and output queues.

Default scheduling is port scheduling. If you configure scheduling instead of using default scheduling, you can configure port scheduling or enhanced transmission selection (ETS) hierarchical port scheduling.

Default scheduling uses weighted round-robin (WRR) scheduling. Each queue receives a portion (weight) of the total available port bandwidth. The scheduling weight is based on the transmit rate (minimum guaranteed bandwidth) of the default scheduler for that queue. For example, queue 7 receives a default scheduling weight of 15 percent of available port bandwidth, and queue 4 receives a default scheduling weight of 35 percent of 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.

You should explicitly map traffic to non-default (unconfigured) queues and schedule bandwidth resources for those queues if you want to use them to forward traffic. By default, queues 1, 2, 5, and 6 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.

If you map traffic to an unconfigured queue and do not schedule bandwidth for the queue, the queue receives only the amount of bandwidth proportional to its default weight (1). The actual amount of bandwidth an unconfigured queue receives 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. Because of their scheduling weights, configured queues have higher priority for bandwidth than unconfigured queues. If a configured queue needs more bandwidth, then less bandwidth is available for unconfigured queues. However, 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 and map it to the forwarding class that is mapped to the queue, and then apply the scheduler map to the port.

Scheduling Priority

Scheduling priority determines the order in which an interface transmits traffic from its output queues. Priority settings ensure that queues containing important traffic receive prioritized access to the outgoing interface bandwidth. The priority setting in the scheduler determines queue priority (a scheduler map maps the scheduler to a forwarding class, the forwarding class is mapped to an output queue, and the output queue uses the CoS properties defined in the scheduler).

By default, all queues are low priority queues. The switch supports three levels of scheduling priority:

- **Low**—In the default CoS state, all queues are low priority queues. Low priority queues transmit traffic based on the weighted round-robin (WRR) algorithm. If you configure scheduling priorities higher than low priority on queues, then the higher priority queues are served before the low priority queues.
- **Medium-low**— (QFX10000 Series switches only) Medium-low priority queues transmit traffic based on the weighted round-robin (WRR) algorithm, and have higher scheduling priority than low priority queues.
- **Medium-high**— (QFX10000 Series switches only) Medium-high priority queues transmit traffic based on the weighted round-robin (WRR) algorithm, and have higher scheduling priority than medium-low priority queues.

- High— (QFX10000 Series switches only) High priority queues transmit traffic based on the weighted round-robin (WRR) algorithm, and have higher scheduling priority than medium-high priority queues.
- Strict-high—You can configure queues as strict-high priority. Strict-high priority queues receive preferential treatment over all other queues, and receive all of their configured bandwidth before other queues are serviced. Other queues do not transmit traffic until strict-high priority queues are empty, and they receive the bandwidth that remains after the strict-high priority queues are serviced. Because strict-high priority queues are always serviced first, strict-high priority queues can starve other queues on a port. Carefully consider how much bandwidth you want to allocate to strict-high priority queues to avoid starving other queues.

NOTE: For QFX10002, QFX10008, and QFX10016 devices, strict-high priority queues share excess bandwidth based on an excess bandwidth sharing weight of 1, which is not configurable. The actual amount of extra bandwidth that strict-high priority traffic exceeding the transmit rate receives depends on how many other queues consume excess bandwidth and the excess rates of those queues.

For QFX10002-60C, excess traffic on the strict-high queue will starve other high/low priority queues.

When you define scheduling priorities for queues instead of using the default priorities (by default all queues are low priority), the switch uses the priorities to determine the order of packet transmission from the queues. The switch services traffic of different scheduling priorities in a strict order, using round-robin (RR) scheduling to arbitrate queue transmission service among queues of the same priority. The switch transmits packets in the following order:

1. Strict-high priority traffic within the configured queue transmit rate (on strict-high priority queues, the transmit rate limits the amount of traffic treated as strict-high priority traffic). When traffic arrives on a strict-high priority queue, the switch forwards it before servicing other queues.
2. High priority traffic within the configured queue transmit rate (on high priority queues, the transmit rate sets the minimum guaranteed bandwidth)
3. Medium-high priority traffic within the configured queue transmit rate (on medium-high priority queues, the transmit rate sets the minimum guaranteed bandwidth)
4. Medium-low priority traffic within the configured queue transmit rate (on medium-low priority queues, the transmit rate sets the minimum guaranteed bandwidth)
5. Low priority traffic within the configured queue transmit rate (on low priority queues, the transmit rate sets the minimum guaranteed bandwidth)
6. All traffic that exceeds the queue transmit rate using weighted round-robin (WRR) scheduling. Traffic that exceeds the queue transmit rate contends for excess port bandwidth (bandwidth that is not consumed after the port meets all guaranteed bandwidth requirements). The switch allocates and

weights excess bandwidth for low priority queues based on the configured queue excess rate, or on the transmit rate if no excess rate is configured. The switch allocates and weights excess bandwidth for strict-high priority queues based on the hard-coded weight “1”, which is not configurable. The actual amount of extra bandwidth that traffic exceeding the transmit rate gets depends on how many other queues consume excess bandwidth and the weighting of those queues.

NOTE: If you use the default CoS configuration, all queues are low priority queues and transmit traffic based on the weighted round-robin (WRR) algorithm.

Bandwidth Scheduling

A queue scheduler allocates port bandwidth to a queue (the scheduler is mapped to a forwarding class, and the forwarding class is mapped to a queue). The bandwidth profile, which consists of minimum guaranteed bandwidth, maximum bandwidth (queue shaping), and excess bandwidth sharing properties configured in the scheduler, defines the amount of port bandwidth a queue can consume during normal and congested transmission periods.

The scheduler regularly reevaluates whether each individual queue is within its defined bandwidth profile by comparing the amount of data the queue receives to the amount of bandwidth the scheduler allocates to the queue. When the received amount is less than the guaranteed minimum amount of bandwidth, the queue is considered to be in profile. A queue is out of profile when its received amount is larger than its guaranteed minimum amount. Out of profile queue data is transmitted only if extra (excess) bandwidth is available. Otherwise, it is buffered if buffer space is available. If no buffer space is available, the traffic might be dropped.

The switch provides features that enable you to control the allocation of port bandwidth to queues, so that you can meet the demands of different types of traffic on a port:

Minimum Guaranteed Bandwidth

The transmit rate determines the minimum guaranteed bandwidth for each forwarding class that is mapped to an output queue, and so determines the minimum bandwidth guarantee on that queue.

If you do not want to use the default configuration, you can set the minimum guaranteed bandwidth in several ways, and with several options, using the `[set class-of-service schedulers scheduler-name transmit-rate (rate | percent percentage) <exact>]` statement:

- **Rate**—Set the minimum guaranteed bandwidth as a fixed amount (rate) in bits-per-second of port bandwidth (for example, 2 Gbps or 800 Mbps).
- **Percent**—Set the minimum guaranteed bandwidth as a percentage of port bandwidth (for example, 25 percent).

- **Exact**—(QFX10000 switches only) Shape the queue to the transmit rate so that the transmit rate is the maximum amount of bandwidth a queue can use. The queue cannot share extra port bandwidth if you configure the exact option. Configuring a transmit rate as *exact* is how you set a shaping rate to configure the maximum amount of bandwidth low and high priority queues can consume, and the maximum is the transmit rate. You cannot use the exact option on a strict-high priority queue.

NOTE: On QFX10000 switches, oversubscribing all 8 queues configured with the transmit rate exact (shaping) statement at the [edit class-of-service schedulers *scheduler-name*] hierarchy level might result in less than 100 percent utilization of port bandwidth.

- **Extra bandwidth sharing**—On low and high priority queues, if you configure an excess rate, the excess rate determines the amount of extra port bandwidth a queue can use. If you do not configure an excess rate, the transmit rate determines how much excess (extra) bandwidth a low and high priority queue can share. If you do not configure an excess rate, then each queue shares extra bandwidth in proportion to its transmit rate.

You cannot configure an excess rate on strict-high priority queues. Strict-high priority queues share extra bandwidth based on a scheduling weight of “1”, which is not configurable. The actual amount of extra bandwidth that traffic exceeding the transmit rate gets depends on how many other queues consume excess bandwidth and the excess rates of those queues.

NOTE: The sum of the transmit rates of the queues on a port should not exceed the total bandwidth of that port. (You cannot guarantee a combined minimum bandwidth for the queues on a port that is greater than the total port bandwidth.)

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 might round small fixed rates to a lower percentage. For example, a fixed rate of 350 Mbps is rounded down to 3 percent.

The bandwidth a low or high priority queue consumes can exceed the configured minimum rate if additional bandwidth is available, and if you do not configure the transmit rate as *exact* on QFX10000 switches. During periods of congestion, the configured transmit rate is the guaranteed minimum bandwidth for the queue. This behavior enables you to ensure that each queue receives the amount of bandwidth appropriate to its required level of service and is also able to share unused bandwidth.

Maximum Bandwidth (Rate Shaping on Low and High Priority Queues and LAGs)

On QFX10000 switches, the optional `exact` keyword in the `[set class-of-service schedulers scheduler-name transmit-rate (rate | percent percentage) <exact>]` configuration statement shapes the transmission rate of low and high priority queues. When you specify the `exact` option, the switch drops traffic that exceeds the configured transmit rate, even if excess bandwidth is available. Rate shaping prevents a queue from using more bandwidth than is appropriate for the planned service level of the traffic on the queue. You cannot use the `exact` option on a strict-high priority queue.

Configuring rate shaping 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.

LAG interfaces consist 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 a 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 rate shaping to LAG traffic, let's look at a LAG interface named `ae0` that has two member links, `xe-0/0/20` and `xe-0/0/21`. On LAG `ae0`, we configure rate shaping of 2g by including the `transmit-rate 2g exact` statement in the queue scheduler, and apply the scheduler to 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. The result is that best-effort traffic receives a maximum of 4g on the LAG, not the 2g set by the rate shaping statement. When hashing spreads the traffic assigned to an output queue (which is mapped to a forwarding class) across multiple LAG links, the effective shaping rate (cumulative maximum bandwidth) on the LAG is:

(number of LAG member interfaces) x (shaping rate for the output queue) = cumulative LAG shaping rate

Limiting Bandwidth Consumed by Strict-High Priority Queues

You can limit the amount of traffic that receives strict-high priority treatment on a queue by configuring a transmit rate on the strict-high priority queue. The transmit rate sets the amount of traffic that receives strict-high priority treatment. 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 gets

depends on how many other queues consume excess bandwidth and the excess rates of those queues. 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: Configuring a transmit rate on a low or high priority queue sets the guaranteed minimum bandwidth of the queue, as described in ["Minimum Guaranteed Bandwidth" on page 375](#).



CAUTION: If you configure strict-high priority queues, we strongly recommend that you configure a transmit rate on the queues to prevent them from starving low and high priority queues on that port. This is especially important if you configure more than one strict-high priority queue on a port. Although it is not mandatory to configure a transmit rate on strict-high priority queues, if you do not configure a transmit rate, the strict-high priority queues can consume all of the port bandwidth and starve the other queues.

Sharing Extra Bandwidth (Excess Rate on Low and High Priority Queues)

Extra bandwidth is essentially the bandwidth remaining after the switch meets all guaranteed bandwidth requirements. Extra bandwidth is available to low and high priority traffic when the queues on a port do not use all of the available port bandwidth.

By default, extra port bandwidth is shared among the forwarding classes on a port in proportion to the transmit rate of each queue. You can explicitly configure the amount of extra bandwidth a queue can share by setting an excess-rate in the scheduler of a low or high priority queue. The configured excess rate overrides the transmit rate and determines the percentage of extra bandwidth the queue can consume.

NOTE: You cannot configure an excess rate on a strict-high priority queue. Strict-high priority queues share excess bandwidth based on an excess bandwidth sharing weight of "1", which is not configurable. The actual amount of extra bandwidth that strict-high priority traffic exceeding the transmit rate receives depends on how many other queues consume excess bandwidth and the excess rates of those queues.

NOTE: QFX 10002, QFX 10008, and QFX 10016 support multiple strict-high queues. QFX 10002-60C supports only one strict-high queue.

An example of extra bandwidth allocation based on transmit rates is a port that has traffic running on three forwarding classes, `best-effort`, `fcoe`, and `network-control`. In this example, the `best-effort` forwarding class has a transmit rate of 2 Gbps, forwarding class `fcoe` has a transmit rate of 4 Gbps, and `network-control` has a transmit rate of 2 Gbps, for a total of 8 Gbps of the port bandwidth. After servicing the minimum guaranteed bandwidth of these three queues, the port has 2 Gbps of available extra bandwidth.

If all three queues still have packets to forward, the queues receive the extra bandwidth in proportion to their transmit rates, so the `best-effort` queue receives an extra 500 Mbps, the `fcoe` queue receives an extra 1 Gbps, and the `network-control` queue receives an extra 500 Mbps.

If you configure an excess rate for a queue, the excess rate determines the proportion of extra bandwidth that the queue receives in the same way that the default (transmit rate) determines the proportion of extra bandwidth a queue receives. In the previous example, if you configured an excess rate of 20 percent on the `fcoe` forwarding class, and the transmit rates of the `best-effort` and `network-control` forwarding classes remained 2g (with no configured excess rate, so the 2g transmit rate for each queue still determines the excess rate), then the 2 Gbps of extra bandwidth would be allocated evenly among the three queues because all three queues have the same excess rate.

In the previous example, if you configured an excess rate of 10 percent on the `fcoe` forwarding class, and the transmit rates of the `best-effort` and `network-control` forwarding classes remained 2g (again with no configured excess rate, so the 2g transmit rate for each queue still determines the excess rate), the 2 Gbps of extra bandwidth would be allocated 800 Mbps to the `best-effort` queue, 400 Mbps to the `fcoe` queue, and 800 Mbps to the `network-control` queue (again, in proportion to the queue excess rates).

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 then 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.

You associate a scheduler map with an interface to apply the drop profiles and other scheduler elements to traffic in the forwarding class mapped to the scheduler on that interface.

Buffer Size

On QFX10000 switches, the buffer size is the amount of time in milliseconds of port bandwidth that a queue can use to continue to transmit packets during periods of congestion, before the buffer runs out and packets begin to drop.

The switch can use up to 100 ms total (combined) buffer space for all queues on a port. A buffer-size configured as one percent is equal to 1 ms of buffer usage. A buffer-size of 15 percent (the default value for the best effort and network control queues) is equal to 15 ms of buffer usage.

The total buffer size of the switch is 4 GB. A 40-Gigabit port can use up to 500 MB of buffer space, which is equivalent to 100 ms of port bandwidth on a 40-Gigabit port. A 10-Gigabit port can use up to 125 MB of buffer space, which is equivalent to 100 ms of port bandwidth on a 10-Gigabit port. The total buffer sizes of the eight output queues on a port cannot exceed 100 percent, which is equal to the full 100 ms total buffer available to a port. The maximum amount of buffer space any queue can use is also 100 ms (which equates to a 100 percent buffer-size configuration), but if one queue uses all of the buffer, then no other queue receives buffer space.

There is no minimum buffer allocation, so you can set the buffer-size to zero (0) for a queue. However, we recommend that on queues on which you enable PFC to support lossless transport, you allocate a minimum of 5 ms (a minimum buffer-size of 5 percent). The two default lossless queues, fcoe and no-loss, have default buffer-size values of 35 ms (35 percent).

NOTE: If you do not configure buffer-size and you do not explicitly configure a queue scheduler, the default buffer-size is the default transmit rate of the queue. If you explicitly configure a queue scheduler, the default buffer allocations are not used. If you explicitly configure a queue scheduler, configure the buffer-size for each queue in the scheduler, keeping in mind that the total buffer-size of the queues cannot exceed 100 percent (100 ms).

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.

Queue buffer allocation is dynamic, shared among ports as needed. However, a queue cannot use more than its configured amount of buffer space. For example, if you are using the default CoS configuration, the best-effort queue receives a maximum of 15 ms of buffer space because the default transmit rate for the best-effort queue is 15 percent.

If a switch experiences congestion, queues continue to receive their full buffer allocation until 90 percent of the 4 GB buffer space is consumed. When 90 percent of the buffer space is in use, the amount of buffer space per port, per queue, is reduced in proportion to the configured buffer size for each queue. As the percentage of consumed buffer space rises above 90 percent, the amount of buffer space per port, per queue, continues to be reduced.

On 40-Gigabit ports, because the total buffer is 4 GB and the maximum buffer a port can use is 500 MB, up to seven 40-Gigabit ports can consume their full 100 ms allocation of buffer space. However, if an eighth 40-Gigabit port requires the full 500 MB of buffer space, then the buffer allocations are proportionally reduced because the buffer consumption is above 90 percent.

On 10-Gigabit ports, because the total buffer is 4 GB and the maximum buffer a port can use is 125 MB, up to 28 10-Gigabit ports can consume their full 100 ms allocation of buffer space. However, if a 29th 10-Gigabit port requires the full 125 MB of buffer space, then the buffer allocations are proportionally reduced because the buffer consumption is above 90 percent.

Explicit Congestion Notification

ECN enables end-to-end congestion notification between two endpoints on TCP/IP based networks. The two endpoints are an ECN-enabled sender and an ECN-enabled receiver. 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.

ECN is disabled by default. Normally, you enable ECN only on queues that handle best-effort traffic because other traffic types use different methods of congestion notification—lossless traffic uses

priority-based flow control (PFC) and strict-high priority traffic receives all of the port bandwidth it requires up to the point of a configured rate (see ["Scheduling Priority" on page 373](#)).

Scheduler Maps

A scheduler map maps a forwarding class to a queue scheduler. After configuring a scheduler, you must include it in a scheduler map, and apply the scheduler map to an interface to implement the configured queue scheduling.

RELATED DOCUMENTATION

[Understanding Junos CoS Components | 21](#)

[Understanding CoS Priority Group Scheduling | 403](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

[Understanding CoS Virtual Output Queues \(VOQs\) | 406](#)

Understanding CoS Explicit Congestion Notification

[Understanding CoS Scheduling Behavior and Configuration Considerations | 332](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)

[Example: Configuring Maximum Output Bandwidth | 431](#)

[Example: Configuring Queue Scheduling Priority | 360](#)

[Example: Configuring Queue Schedulers for Port Scheduling | 386](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring WRED Drop Profiles | 286](#)

Example: Configuring ECN

Defining CoS Queue Schedulers for Port Scheduling

Schedulers define the CoS properties of output queues. You configure CoS properties in a scheduler, then map the scheduler to a forwarding class. Forwarding classes are in turn mapped to output queues. Classifiers map incoming traffic into forwarding classes based on IEEE 802.1p, DSCP, or EXP code points. 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:

- **priority**—One of three bandwidth priorities that queues associated with a scheduler can receive:
 - **low**—The scheduler has low priority.
 - **high**—The scheduler has high priority. High priority traffic takes precedence over low priority traffic.
 - **strict-high**—The scheduler has strict-high priority. Strict-high priority queues receive preferential treatment over low-priority queues and receive all of their configured bandwidth before low-priority queues are serviced. Low-priority queues do not transmit traffic until strict-high priority queues are empty.

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.

- **transmit-rate**—Minimum guaranteed bandwidth, also known as the *committed information rate (CIR)*, set as a percentage rate or as an absolute value in bits per second. By default, the transmit rate also determines the amount of excess (extra) port bandwidth the queue can share if you do not explicitly configure an excess rate. Extra bandwidth is allocated among the queues on the port in proportion to the transmit rate of each queue. Except on QFX10000 switches, you can configure *shaping-rate* to throttle the rate of packet transmission. On QFX10000 switches, on queues that are not strict-high priority queues, you can configure a transmit rate as *exact*, which shapes the transmission by setting the transmit rate as the maximum bandwidth the queue can consume on the port.

NOTE: On QFX10000 switches, oversubscribing all 8 queues configured with the transmit rate exact (shaping) statement at the [edit class-of-service schedulers *scheduler-name*] hierarchy level might result in less than 100 percent utilization of port bandwidth.

On strict-high priority queues, the transmit rate sets the amount of bandwidth used for strict-high priority forwarding; traffic in excess of the transmit rate is treated as best-effort traffic that receives the queue excess rate.

NOTE: Include the preamble bytes and interframe gap (IFG) bytes as well as the data bytes in your bandwidth calculations.

- **excess-rate**—Percentage of extra bandwidth (bandwidth that is not used by other queues) a low-priority queue can receive. If not set, the switch uses the transmit rate to determine extra bandwidth sharing. You cannot set an excess rate on a strict-high priority queue.
- **drop-profile-map**—Drop profile mapping to a packet loss priority 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**—ECN enable 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: 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.

To apply scheduling properties to traffic, map schedulers to forwarding classes using a scheduler map, and then apply the scheduler map to interfaces. Using different scheduler maps, you can map different schedulers to 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; optionally, set a maximum bandwidth limit (shaping rate) on a low priority queue by configuring either *shaping-rate* (except on QFX10000 switches) or the *exact* option (only on QFX10000 switches):

```
[edit class-of-service]
user@switch# set schedulers scheduler-name transmit-rate (rate | percent percentage)
<exact>
```

2. Set the amount of excess bandwidth a low-priority queue can share:

```
[edit class-of-service]
user@switch# set schedulers scheduler-name excess-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) drop-profile drop-profile-name
```

5. Configure the size of the 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.

```
[edit class-of-service]
user@switch# set interfaces interface-name scheduler-map scheduler-map-name
```

RELATED DOCUMENTATION

[Example: Configuring Queue Schedulers for Port Scheduling | 386](#)

Example: Configuring ECN

[Defining CoS Queue Scheduling Priority | 358](#)

[Configuring CoS WRED Drop Profiles | 284](#)

[Monitoring CoS Scheduler Maps | 365](#)

[Understanding CoS Port Schedulers | 368](#)

Understanding CoS Explicit Congestion Notification

Example: Configuring Queue Schedulers for Port Scheduling

IN THIS SECTION

- [Requirements | 388](#)
- [Overview | 388](#)
- [Verification | 391](#)

Schedulers define the CoS properties of output queues. You configure CoS properties in a scheduler, then map the scheduler to a forwarding class. Forwarding classes are in turn mapped to output queues. Classifiers map incoming traffic into forwarding classes based on IEEE 802.1p, DSCP, or EXP code points. 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 buffer-size percent 20
set schedulers be-sched excess-rate 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 interfaces xe-0/0/7 scheduler-map be-map
```

Step-by-Step Procedure

To configure a CoS scheduler:

1. Create scheduler (be-sched) 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 buffer-size percent 20
user@switch# set be-sched excess-rate 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 with an interface to apply scheduling to the best-effort forwarding class output queue:

```
[edit class-of-service]
set interfaces xe-0/0/7 scheduler-map be-map
```

Requirements

This example uses the following hardware and software components:

- One QFX10000 switch.
- Junos OS Release 15.1X53-D10 or later for the QFX Series

Overview

Scheduler parameters define the following characteristics for the queues mapped to the scheduler:

The parameters you configure in a scheduler define the following characteristics for the queues mapped to the scheduler:

- **priority**—One of three bandwidth priorities that queues associated with a scheduler can receive:
 - **low**—The scheduler has low priority.
 - **high**—The scheduler has high priority. High priority traffic takes precedence over low priority traffic.
 - **strict-high**—The scheduler has strict-high priority. Strict-high priority queues receive preferential treatment over low-priority queues and receive all of their configured bandwidth before low-priority queues are serviced. Low-priority queues do not transmit traffic until strict-high priority queues are empty.

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.

- **transmit-rate**—Minimum guaranteed bandwidth, also known as the *committed information rate (CIR)*, set as a percentage rate or as an absolute value in bits per second. By default, the transmit rate also determines the amount of excess (extra) port bandwidth the queue can share if you do not explicitly configure an excess rate. Extra bandwidth is allocated among the queues on the port in proportion to the transmit rate of each queue. On queues that are not strict-high priority queues, you can configure a transmit rate as exact, which shapes the transmission by setting the transmit rate as the maximum bandwidth the queue can consume on the port.

On strict-high priority queues, the transmit rate sets the amount of bandwidth used for strict-high priority forwarding; traffic in excess of the transmit rate is treated as best-effort traffic that receives the queue excess rate.

NOTE: Include the preamble bytes and interframe gap (IFG) bytes as well as the data bytes in your bandwidth calculations.

- **excess-rate**—Percentage of extra bandwidth (bandwidth that is not used by other queues) a low-priority queue can receive. If not set, the switch uses the transmit rate to determine extra bandwidth sharing. You cannot set an excess rate on a strict-high priority queue.
- **drop-profile-map**—Drop profile mapping to a packet loss priority 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—ECN enable 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: 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.

Scheduler maps map schedulers to forwarding classes, and forwarding classes are mapped to output queues. After you configure schedulers and map them to forwarding classes in a scheduler map, you attach the scheduler map to an interface to implement the configured scheduling on output queues on that interface.

This process configures the bandwidth properties, scheduling, priority, and WRED characteristics that you map to forwarding classes (and thus to output queues) in a scheduler map.

[Table 78 on page 390](#) shows the configuration components for this example.

Table 78: Components of the Port Output Queue Scheduler Configuration Example

Component	Settings
Hardware	One switch
Scheduler	Name: be-sched Transmit rate: 20% Buffer size: 20% Excess rate: 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

Verification

IN THIS SECTION

- [Verifying the Scheduler Configuration | 391](#)
- [Verifying the Scheduler Map Configuration | 391](#)
- [Verifying That the Scheduler Is Associated with the Interface | 392](#)

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 (`transmit-rate`) of 2 Gbps, an extra bandwidth sharing rate (`excess-rate`) of 20 percent, 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;
buffer-size percent 20;
excess-rate 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`.

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;
```

Verifying That the Scheduler Is Associated with the Interface

Purpose

Verify that the scheduler map `be-sched` is attached to 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
scheduler-map be-map;
```

RELATED DOCUMENTATION

[Example: Configuring WRED Drop Profiles | 286](#)

Example: Configuring ECN

[Defining CoS Queue Schedulers for Port Scheduling | 382](#)

[Monitoring CoS Scheduler Maps | 365](#)

[Understanding CoS Port Schedulers | 368](#)

[Understanding CoS Virtual Output Queues \(VOQs\) | 406](#)

CoS Port Shaping

SUMMARY

Port shaping enables you to control the amount of traffic passing through an interface. Port shaping enables you to shape the aggregate traffic through an interface to a rate that is less than the line rate for that interface. This can be useful to reduce downstream congestion.

IN THIS SECTION

- [Understanding Port Shaping | 393](#)
- [Configuring Port Shaping | 393](#)

This topic describes port shaping and how to configure port shaping.

Understanding Port Shaping

By default, shaping is not configured on an interface and traffic can be up to the line rate for that interface. When you configure port shaping on an interface, you specify a value that indicates the maximum amount of traffic that can pass through the interface.

Port shaping enables you to shape the aggregate traffic through a port or channel to a rate that is less than the line rate. You specify the port shaping rate as the peak rate at which traffic can pass through the interface. You specify the rate as a value in bits per second (bps) either as a decimal number or as a decimal number followed by the abbreviation k (1000), m (1,000,000), or g (1,000,000,000) and the value can range from 1000 through 160,000,000,000 bps. This value must be less than the maximum bandwidth for that interface.

You can configure port shaping on network interfaces, aggregated Ethernet interfaces (also known as link aggregation groups (LAGs)), and loopback interfaces.

NOTE: On EX4650, QFX5110, QFX5120, QFX5200, QFX5210 Series switches, when you configure a shaping rate on an aggregated Ethernet (ae) interface, all members of the ae interface are shaped at the configured shaping rate. For example, consider an interface ae0 that consists of three interfaces: xe-0/0/0, xe-0/0/1, and xe-0/0/2. If you configure a shaping rate of X Mbps on ae0, traffic up to the rate of X Mbps flows through each of the three interfaces. Therefore, the total traffic flowing through ae0 can be at the rate of 3X Mbps.

Configuring Port Shaping

You can configure port shaping on network interfaces, aggregated Ethernet interfaces (also known as link aggregation groups (LAGs)), and loopback interfaces.

To configure port shaping on an interface:

1. Ensure that the interface on which you want to configure port shaping is up and running.
2. Assign a `shaping-rate` for the interface:

```
[edit class-of-service]
user@switch# set interfaces interface-name shaping-rate value
```

For example:

```
[edit class-of-service]
user@switch# set interfaces xe-0/0/0 shaping-rate 3200000000000
```

The value indicates the maximum amount of traffic (in bps) that can pass through the interface. This value must be less than the maximum bandwidth for that interface.

3. Commit your changes.
4. Verify your configuration.

For example:

```
[edit class-of-service]
user@host# show
interfaces {
  xe-0/0/0 {
    shaping-rate 3200g;
  }
}
```

5. Run `show class-of-service interface interface-name` to verify the shaping rate for the interface.

For example:

```
user@host> show class-of-service interface xe-0/0/0
Physical interface: xe-0/0/0, Index: 650
Maximum usable queues: 10, Queues in use: 5
Exclude aggregate overhead bytes: disabled
Logical interface aggregate statistics: disabled
Shaping rate: 3200000000000 bps
Scheduler map: <default>, Index: 2
Congestion-notification: Disabled
```

Logical interface: xe-0/0/0.16386, Index: 874

RELATED DOCUMENTATION

shaping-rate (Applying to an Interface)

show class-of-service interface

Troubleshooting Egress Bandwidth Issues

IN THIS CHAPTER

- [Troubleshooting Egress Bandwidth That Exceeds the Configured Minimum Bandwidth | 396](#)
- [Troubleshooting Egress Bandwidth That Exceeds the Configured Maximum Bandwidth | 398](#)
- [Troubleshooting Egress Queue Bandwidth Impacted by Congestion | 399](#)

Troubleshooting Egress Bandwidth That Exceeds the Configured Minimum Bandwidth

IN THIS SECTION

- [Problem | 396](#)
- [Cause | 397](#)
- [Solution | 397](#)

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](#)

Troubleshooting Egress Bandwidth That Exceeds the Configured Maximum Bandwidth

IN THIS SECTION

● [Problem | 398](#)

● [Cause | 398](#)

● [Solution | 399](#)

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 | 399](#)
- [Cause | 399](#)
- [Solution | 400](#)

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](#)

Traffic Control Profiles and Priority Group Scheduling

IN THIS CHAPTER

- [Understanding CoS Traffic Control Profiles | 401](#)
- [Understanding CoS Priority Group Scheduling | 403](#)
- [Understanding CoS Virtual Output Queues \(VOQs\) | 406](#)
- [Defining CoS Traffic Control Profiles \(Priority Group Scheduling\) | 412](#)
- [Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)
- [Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth | 417](#)
- [Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)
- [Understanding CoS Priority Group Shaping and Queue Shaping \(Maximum Bandwidth\) | 428](#)
- [Example: Configuring Maximum Output Bandwidth | 431](#)

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 \(ETS\) | 438](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Defining CoS Traffic Control Profiles \(Priority Group Scheduling\) | 412](#)

Understanding CoS Priority Group Scheduling

IN THIS SECTION

- [Priority Group Scheduling Components | 404](#)
- [Default Traffic Control Profile | 404](#)
- [Guaranteed Rate \(Minimum Guaranteed Bandwidth\) | 405](#)
- [Sharing Extra Bandwidth | 405](#)
- [Shaping Rate \(Maximum Bandwidth\) | 405](#)
- [Scheduler Maps | 406](#)

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 79 on page 404](#) provides a quick reference to the traffic control profile components you can configure to determine the bandwidth properties of priority groups, and [Table 80 on page 404](#) provides a quick reference to some related scheduling configuration components.

Table 79: 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 80: 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 | 21](#)

[Understanding CoS Output Queue Schedulers | 338](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

[Understanding CoS Scheduling Behavior and Configuration Considerations | 332](#)

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)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)

[Example: Configuring Maximum Output Bandwidth | 431](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring WRED Drop Profiles | 286](#)

[Example: Configuring Drop Profile Maps | 293](#)

Understanding CoS Virtual Output Queues (VOQs)

IN THIS SECTION

● [VOQ Architecture | 407](#)

The traditional method of forwarding traffic through a switch is based on buffering ingress traffic in input queues on ingress interfaces, forwarding the traffic across the switch fabric to output queues on egress interfaces, and then buffering traffic again on the output queues before transmitting the traffic to the next hop. The traditional method of queueing packets on an ingress port is storing traffic destined for different egress ports in the same input queue (buffer).

During periods of congestion, the switch might drop packets at the egress port, so the switch might spend resources transporting traffic across the switch fabric to an egress port, only to drop that traffic instead of forwarding it. And because input queues store traffic destined for different egress ports, congestion on one egress port could affect traffic on a different egress port, a condition called *head-of-line blocking (HOLB)*.

Virtual output queue (VOQ) architecture takes a different approach:

- Instead of separate physical buffers for input and output queues, the switch uses the physical buffers on the ingress pipeline of each Packet Forwarding Engine (PFE) chip to store traffic for every egress port. Every output queue on an egress port has buffer storage space on every ingress pipeline on all of the PFE chips on the switch. The mapping of ingress pipeline storage space to output queues is 1-to-1, so each output queue receives buffer space on each ingress pipeline.
- Instead of one input queue containing traffic destined for multiple different output queues (a one-to-many mapping), each output queue has a dedicated VOQ comprised of the input buffers on each packet forwarding chip that are dedicated to that output queue (a 1-to-1 mapping). This architecture prevents communication between any two ports from affecting another port.
- Instead of storing traffic on a physical output queue until it can be forwarded, a VOQ does not transmit traffic from the ingress port across the fabric to the egress port until the egress port has the resources to forward the traffic.

A VOQ is a collection of input queues (buffers) that receive and store traffic destined for one output queue on one egress port. Each output queue on each egress port has its own dedicated VOQ, which consists of all of the input queues that are sending traffic to that output queue.

VOQ Architecture

A VOQ represents the ingress buffering for a particular output queue. A unique buffer ID identifies each output queue on a PFE chip. Each of the six PFE chips uses the same unique buffer ID for a particular output queue. The traffic stored using a particular buffer ID on the six PFE chips comprises the traffic destined for one particular output queue on one port, and is the VOQ for that output queue.

A switch that has 72 egress ports with 8 output queues on each port, has 576 VOQs on each PFE chip ($72 \times 8 = 576$). Because the switch has six PFE chips, the switch has a total of 3,456 VOQs ($576 \times 6 = 3,456$).

A VOQ is distributed across all of the PFE chips that are actively sending traffic to that output queue. Each output queue is the sum of the total buffers assigned to that output queue (by its unique buffer ID) across all of the PFE chips. So the output queue itself is virtual, not physical, although the output queue is comprised of physical input queues.

Round-Trip Time Buffering

Although there is no output queue buffering during periods of congestion (no long-term storage), there is a small physical output queue buffer on egress line cards to accommodate the round-trip time for traffic to traverse the switch fabric from ingress to egress. The round-trip time consists of the time it takes the ingress port to request egress port resources, receive a grant from the egress port for resources, and transmit the data across the switch fabric.

That means if a packet is not dropped at the switch ingress, and the switch forwards the packet across the fabric to the egress port, the packet will not be dropped and will be forwarded to the next hop. All packet drops take place in the ingress pipeline.

The switch has 4 GB of external DRAM to use as a delay bandwidth buffer (DBB). The DBB provides storage for ingress ports until the ports can forward traffic to egress ports.

Requesting and Granting Egress Port Bandwidth

When packets arrive at an ingress port, the ingress pipeline stores the packet in the ingress queue with the unique buffer ID of the destination output queue. The switch makes the buffering decision after performing the packet lookup. If the packet belongs to a class for which the maximum traffic threshold has been exceeded, the packet might not be buffered and might be dropped. To transport packets across the switch fabric to egress ports:

1. The ingress line card PFE request scheduler sends a request to the egress line card PFE grant scheduler to notify the egress PFE that data is available for transmission.
2. When there is available egress bandwidth, the egress line card grant scheduler responds by sending a bandwidth grant to the ingress line card PFE.
3. The ingress line card PFE receives the grant from the egress line card PFE, and transmits the data to the egress line card.

Ingress packets remain in the VOQ on the ingress port input queues until the output queue is ready to accept and forward more traffic.

Under most conditions, the switch fabric is fast enough to be transparent to egress class-of-service (CoS) policies, so the process of forwarding traffic from the ingress pipeline, across the switch fabric, to egress ports, does not affect the configured CoS policies for the traffic. The fabric only affects CoS policy if there is a fabric failure or if there is an issue of port fairness.

When a packet ingresses and egresses the same PFE chip (local switching), the packet does not traverse the switch fabric. However, the switch uses the same request and grant mechanism to receive egress bandwidth as packets that cross the fabric, so locally switched packets and packets that arrive at a PFE chip after crossing the switch fabric are treated fairly when the traffic is contending for the same output queue.

VOQ Advantages

VOQ architecture provides two major advantages:

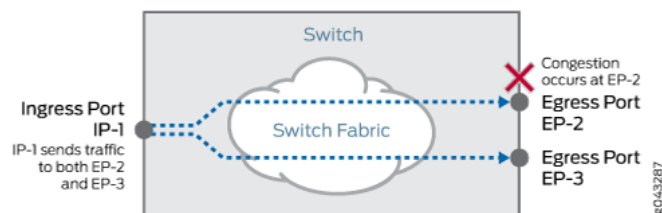
Eliminate Head-of-Line Blocking

VOQ architecture eliminates head-of-line blocking (HOLB) issues. On non-VOQ switches, HOLB occurs when congestion at an egress port affects a different egress port that is not congested. HOLB occurs when the congested port and the uncongested port share the same input queue on an ingress interface.

An example of a HOLB scenario is a switch that has streams of traffic entering one ingress port (IP-1) that are destined for two different egress ports (EP-2 and EP-3):

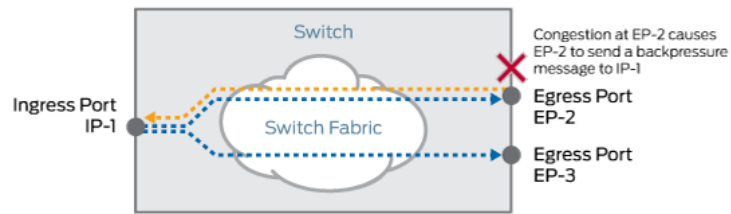
1. Congestion occurs on egress port EP-2. There is no congestion on egress port EP-3, as shown in [Figure 11 on page 409](#).

Figure 11: Congestion Occurs on EP-2



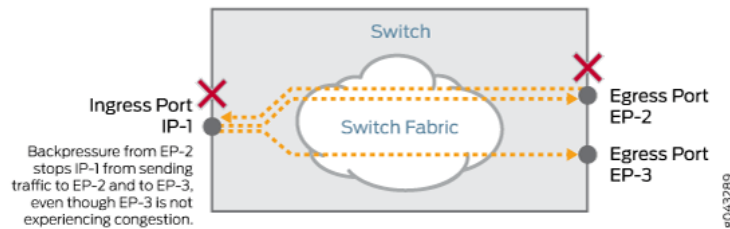
2. Egress port EP-2 sends a backpressure signal to ingress port IP-1, as shown in [Figure 12 on page 410](#).

Figure 12: EP-2 Backpressures IP-1



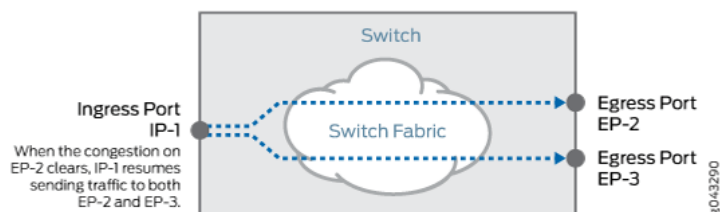
3. The backpressure signal causes the ingress port IP-1 to stop sending traffic and to buffer traffic until it receives a signal to resume sending, as shown in [Figure 13 on page 410](#). Traffic that arrives at ingress port IP-1 destined for uncongested egress port EP-3 is buffered along with the traffic destined for congested port EP-2, instead of being forwarded to port EP-3.

Figure 13: Backpressure from EP-2 Causes IP-1 to Buffer Traffic Instead of Sending Traffic, Affecting EP-3



4. Ingress port IP-1 transmits traffic to uncongested egress port EP-3 only when egress port EP-2 clears enough to allow ingress port IP-1 to resume sending traffic, as shown in [Figure 14 on page 410](#).

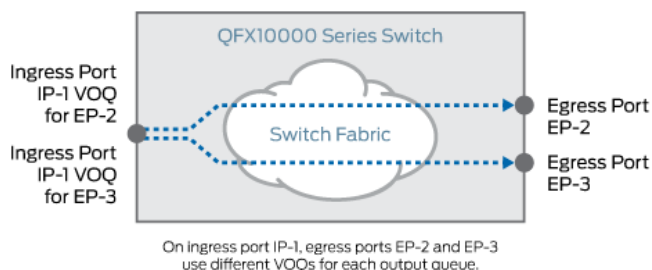
Figure 14: Congestion on EP-2 Clears, Allowing IP-1 to Resume Sending Traffic to Both Egress Ports



In this way, congested egress port EP-2 negatively affects uncongested egress port EP-3, because both egress ports share the same input queue on ingress port IP-1.

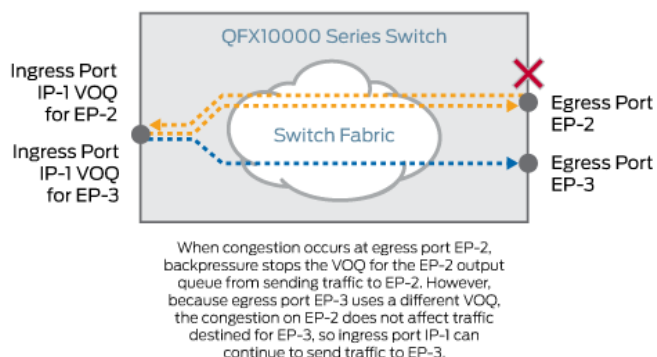
VOQ architecture avoids HOLB by creating a different dedicated virtual queue for each output queue on each interface, as shown in [Figure 15 on page 411](#).

Figure 15: Each Egress Port Has a Separate Virtual Output Queue on IP-1



Because different egress queues do not share the same input queue, a congested egress queue on one port cannot affect an egress queue on a different port, as shown in [Figure 16 on page 411](#). (For the same reason, a congested egress queue on one port cannot affect another egress queue on the same port—each output queue has its own dedicated virtual output queue composed of ingress interface input queues.)

Figure 16: Congestion on EP-2 Does Not Affect Uncongested Port EP-3



Performing queue buffering at the ingress interface ensures that the switch only sends traffic across the fabric to an egress queue if that egress queue is ready to receive that traffic. If the egress queue is not ready to receive traffic, the traffic remains buffered at the ingress interface.

Increase Fabric Efficiency and Utilization

Traditional output queue architecture has some inherent inefficiencies that VOQ architecture addresses.

- **Packet buffering**—Traditional queueing architecture buffers each packet twice in long-term DRAM storage, once at the ingress interface and once at the egress interface. VOQ architecture buffers each packet only once in long-term DRAM storage, at the ingress interface. The switch fabric is fast enough to be transparent to egress CoS policies, so instead of buffering packets a second time at the egress interface, the switch can forward traffic at a rate that does not require deep egress buffers, without affecting the configured egress CoS policies (scheduling).
- **Consumption of resources**—Traditional queueing architecture sends packets from the ingress interface input queue (buffer), across the switch fabric, to the egress interface output queue (buffer). At the egress interface, packets might be dropped, even though the switch has expended resources transporting the packets across the fabric and storing them in the egress queue. VOQ architecture does not send packets across the fabric to the egress interface until the egress interface is ready to transmit the traffic. This increases system utilization because no resources are wasted transporting and storing packets that are dropped later.

Independent of VOQ architecture, the Juniper Networks switching architecture also provides better fabric utilization because the switch converts packets into cells. Cells have a predictable size, which enables the switch to spray the cells evenly across the fabric links and more fully utilize the fabric links. Packets vary greatly in size, and packet size is not predictable. Packet-based fabrics can deliver no better than 65-70 percent utilization because of the variation and unpredictability of packet sizes. Juniper Networks' cell-based fabrics can deliver a fabric utilization rate of almost 95 percent because of the predictability of and control over cell size.

RELATED DOCUMENTATION

[Understanding CoS Port Schedulers | 368](#)

[Example: Configuring Queue Schedulers for Port Scheduling | 386](#)

[Understanding Default CoS Scheduling and Classification | 321](#)

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 \(ETS\) | 445](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)

[Example: Configuring Maximum Output Bandwidth | 431](#)

[Defining CoS Queue Schedulers | 346](#)

[Understanding CoS Traffic Control Profiles | 401](#)

Example: Configuring Traffic Control Profiles (Priority Group Scheduling)

IN THIS SECTION

- [Requirements | 415](#)
- [Overview | 415](#)
- [Verification | 416](#)

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 | 417](#)

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 \(ETS\) | 445](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)

[Example: Configuring Maximum Output Bandwidth | 431](#)

[Example: Configuring Queue Schedulers | 350](#)

[Defining CoS Traffic Control Profiles \(Priority Group Scheduling\) | 412](#)

[Understanding CoS Traffic Control Profiles | 401](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

Understanding CoS Priority Group and Queue Guaranteed Minimum Bandwidth

IN THIS SECTION

● [Guaranteeing Bandwidth Using Hierarchical Scheduling | 418](#)

- [Priority Group Guaranteed Rate \(Guaranteed Minimum Bandwidth\) | 420](#)
- [Queue Transmit Rate \(Guaranteed Minimum Bandwidth\) | 420](#)

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.

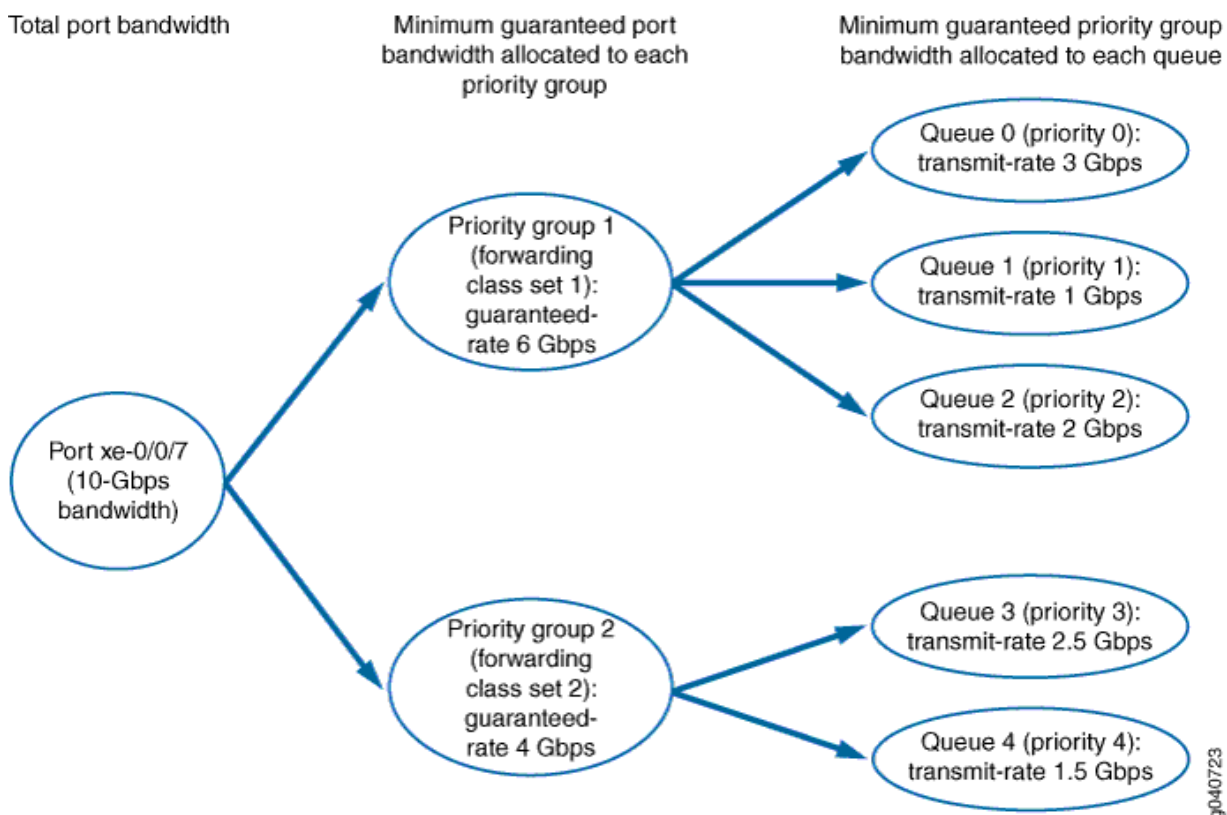
Guaranteeing Bandwidth Using Hierarchical Scheduling

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 17 on page 419](#) 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 17: 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 | 338](#)[Understanding CoS Traffic Control Profiles | 401](#)[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)[Example: Configuring Queue Schedulers | 350](#)[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)[Defining CoS Queue Schedulers | 346](#)[Defining CoS Traffic Control Profiles \(Priority Group Scheduling\) | 412](#)

Example: Configuring Minimum Guaranteed Output Bandwidth

IN THIS SECTION

- [Requirements | 423](#)
- [Overview | 423](#)
- [Verification | 425](#)

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 81 on page 424](#) shows the configuration components for this example:

Table 81: Components of the Minimum Guaranteed Output Bandwidth Configuration Example

Component	Settings
Hardware	QFX3500 switch

Table 81: Components of the Minimum Guaranteed Output Bandwidth Configuration Example
(Continued)

Component	Settings
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 | 426](#)
- [Verifying the Priority Group Minimum Guaranteed Bandwidth and Scheduler Map Association | 426](#)
- [Verifying the Scheduler Map Configuration | 427](#)
- [Verifying Queue \(Forwarding Class\) Membership in the Priority Group | 427](#)
- [Verifying the Egress Interface Configuration | 427](#)

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 \(ETS\) | 445](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring Queue Scheduling Priority | 360](#)

[Example: Configuring Forwarding Class Sets | 184](#)

[Understanding CoS Traffic Control Profiles | 401](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

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

IN THIS SECTION

- [Priority Group Shaping | 429](#)
- [Queue Shaping | 429](#)
- [Shaping Maximum Bandwidth Using Hierarchical Scheduling | 430](#)

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.

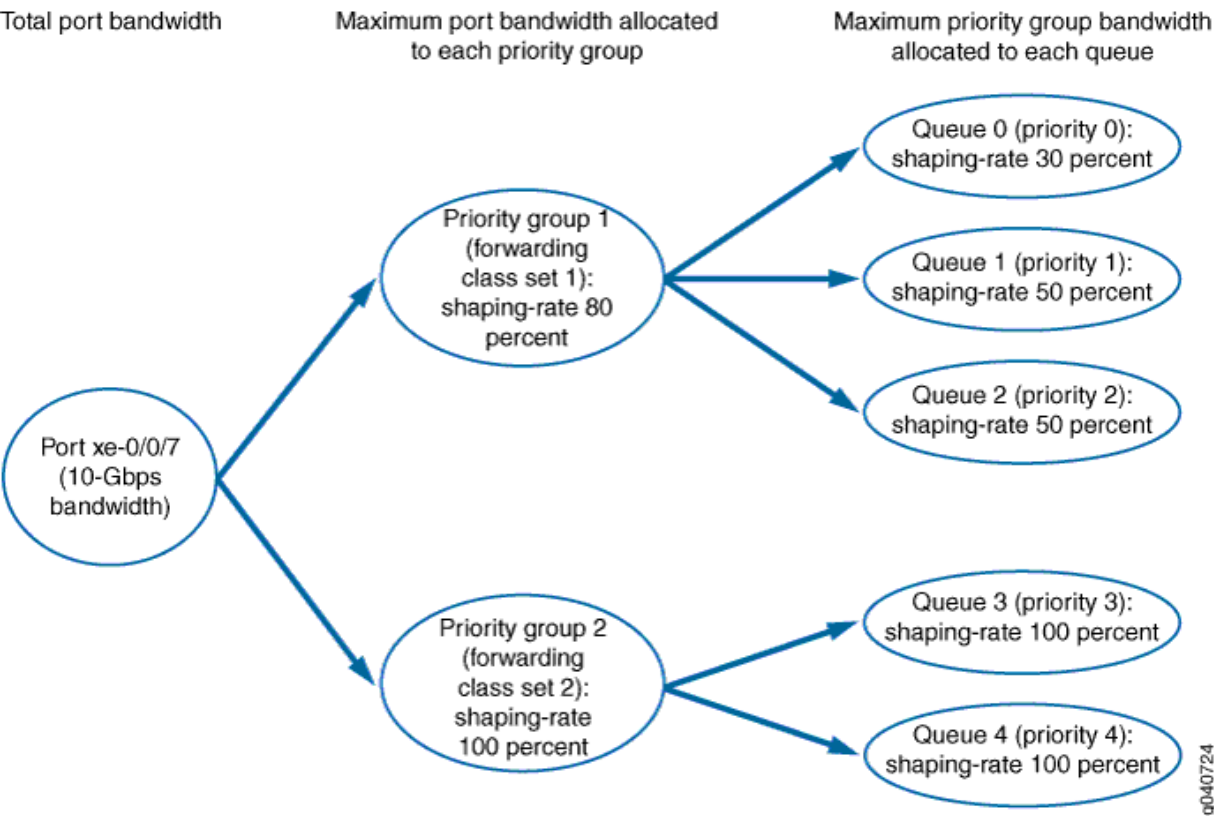
You configure the queue shaping-rate in the scheduler configuration, and you set the shaping-rate for priority groups in the traffic control profile configuration.

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 18 on page 430 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 18: Setting Maximum Bandwidth Using Hierarchical Scheduling



RELATED DOCUMENTATION

[Understanding CoS Output Queue Schedulers | 338](#)

[Understanding CoS Traffic Control Profiles | 401](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Defining CoS Queue Schedulers | 346](#)

[Defining CoS Traffic Control Profiles \(Priority Group Scheduling\) | 412](#)

Example: Configuring Maximum Output Bandwidth

IN THIS SECTION

- [Requirements | 433](#)
- [Overview | 433](#)
- [Verification | 434](#)

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 82 on page 434 shows the configuration components for this example:

Table 82: 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 | 435

- [Verifying the Priority Group Maximum Bandwidth and Scheduler Map Association | 435](#)
- [Verifying the Scheduler Map Configuration | 436](#)
- [Verifying Queue \(Forwarding Class\) Membership in the Priority Group | 436](#)
- [Verifying the Egress Interface Configuration | 437](#)

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 \(ETS\) | 445](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring Forwarding Class Sets | 184](#)

[Understanding CoS Traffic Control Profiles | 401](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

Hierarchical Port Scheduling (ETS)

IN THIS CHAPTER

- Understanding CoS Hierarchical Port Scheduling (ETS) | 438
- Example: Configuring CoS Hierarchical Port Scheduling (ETS) | 445
- Disabling the ETS Recommendation TLV | 480

Understanding CoS Hierarchical Port Scheduling (ETS)

IN THIS SECTION

- Hierarchical Scheduling Tiers | 439
- Hierarchical Scheduling and ETS | 440
- ETS Advertisement in DCBX | 441
- Hierarchical Scheduling Process | 441
- Strict-High Priority Queues and Hierarchical Scheduling | 443
- Default Hierarchical Scheduling | 444

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 83 on page 439](#).

Table 83: 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.

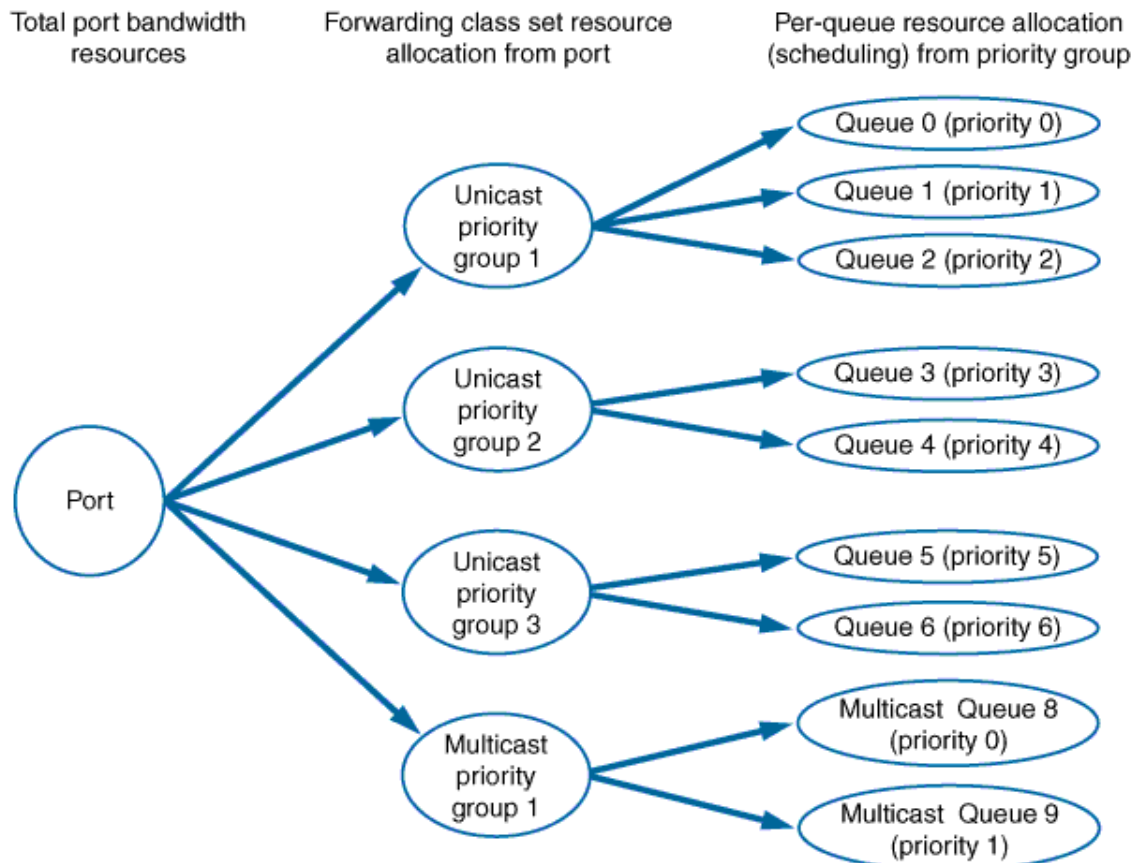
NOTE: If you explicitly configure one or more priority groups on an interface, any priority (forwarding class) that is not assigned to a priority group (forwarding class set) on that interface is assigned to an automatically generated default priority group and receives *no bandwidth*. This means that 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 19 on page 440](#) shows the relationship of port resource allocation to priority groups, and priority group resource allocation to queues (priorities).

Figure 19: 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 20 on page 443](#):

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.

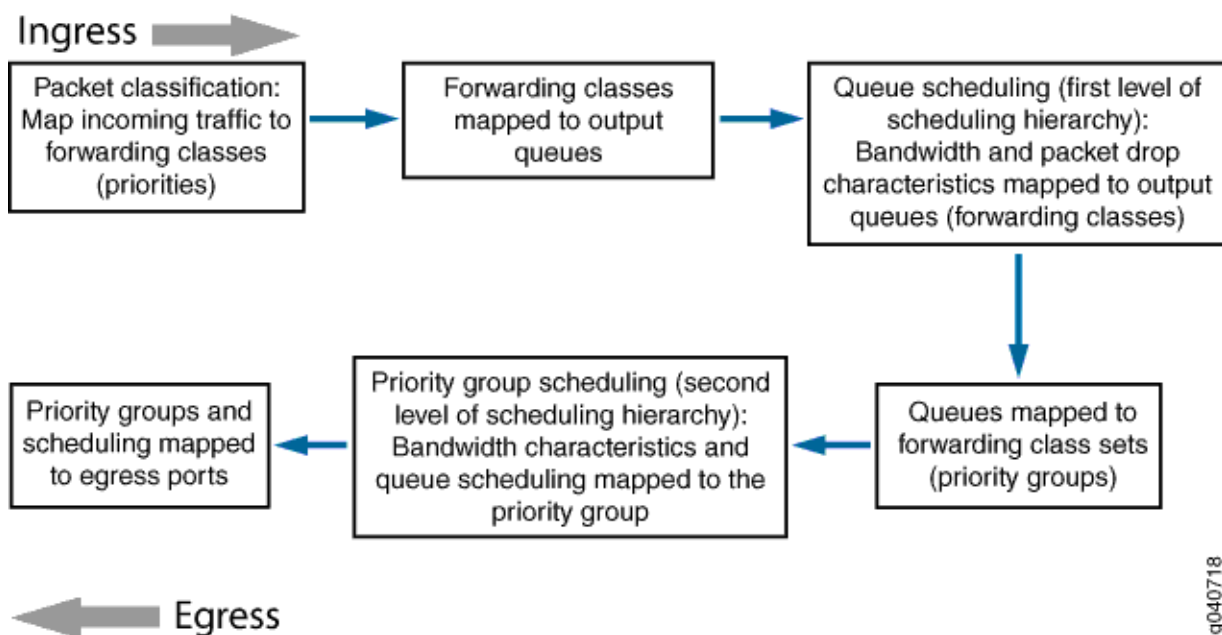
2. Configure the output queues for the forwarding classes (priorities). This consists of either using the default forwarding classes and forwarding-class-to-queue mapping, or creating your own forwarding classes and mapping them to output queues.
3. Allocate resources to the forwarding classes:
 - Define 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 20: 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 multdestination 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 multdestination 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 | 26](#)

[Understanding CoS Output Queue Schedulers | 338](#)

[Understanding CoS Priority Group Scheduling | 403](#)

[Benefits of Configuring CoS Hierarchical Port Scheduling](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

[Understanding CoS Classifiers | 96](#)

[Understanding Default CoS Scheduling and Classification | 321](#)

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)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\) | 445](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)

[Example: Configuring Maximum Output Bandwidth | 431](#)

Example: Configuring CoS Hierarchical Port Scheduling (ETS)

IN THIS SECTION

- [Requirements | 446](#)
- [Overview | 446](#)
- [Configuration | 452](#)
- [Verification | 466](#)

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](#) | 447

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 84 on page 448](#) 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 84: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology

Property	Settings
Hardware	QFX3500 switch
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>

Table 84: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology (Continued)

Property	Settings
Behavior aggregate classifier (maps forwarding classes and loss priorities to incoming packets by IEEE 802.1 code point)	Name—hsclassifier1 Code point mapping: <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
PFC	Congestion notification profile name—gd-cnp PFC enabled on code points: 011 (fcoe priority), 010 (no-loss priority)
Drop profiles NOTE: The fcoe and no-loss priorities (queues) do not use drop profiles because they are lossless traffic classes.	dp-be-low: drop start point 25, drop end point 50, maximum drop rate 80 dp-be-high: drop start point 10, drop end point 40, maximum drop rate 100 dp-hpc: drop start point 75, drop end point 90, maximum drop rate 75 dp-nc: drop start point 80, drop end point 100, maximum drop rate 100
Queue schedulers	be-sched: minimum bandwidth 3g, maximum bandwidth 100%, priority low, drop profiles dp-be-low and dp-be-high fcoe-sched: minimum bandwidth 2.5g, maximum bandwidth 100%, priority low hpc-sched: minimum bandwidth 2g, maximum bandwidth 100%, priority low, drop profile dp-hpc nc-sched: minimum bandwidth 500m, maximum bandwidth 100%, priority low, drop profile dp-nc nl-sched: minimum bandwidth 2g, maximum bandwidth 100%, priority low

Table 84: Components of the Hierarchical Port Scheduling (ETS) Configuration Topology (Continued)

Property	Settings
Forwarding class-to-scheduler mapping	<p>Scheduler map be-map: Forwarding class best-effort, scheduler be-sched Forwarding class be2, scheduler be-sched Forwarding class network-control, scheduler nc-sched</p> <p>Scheduler map gd-map: Forwarding class fcoe, scheduler fcoe-sched Forwarding class no-loss, scheduler nl-sched</p> <p>Scheduler map hpc-map: Forwarding class hpc, scheduler hpc-sched</p>
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 21 on page 451 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 21: Hierarchical Port Scheduling Components Block Diagram

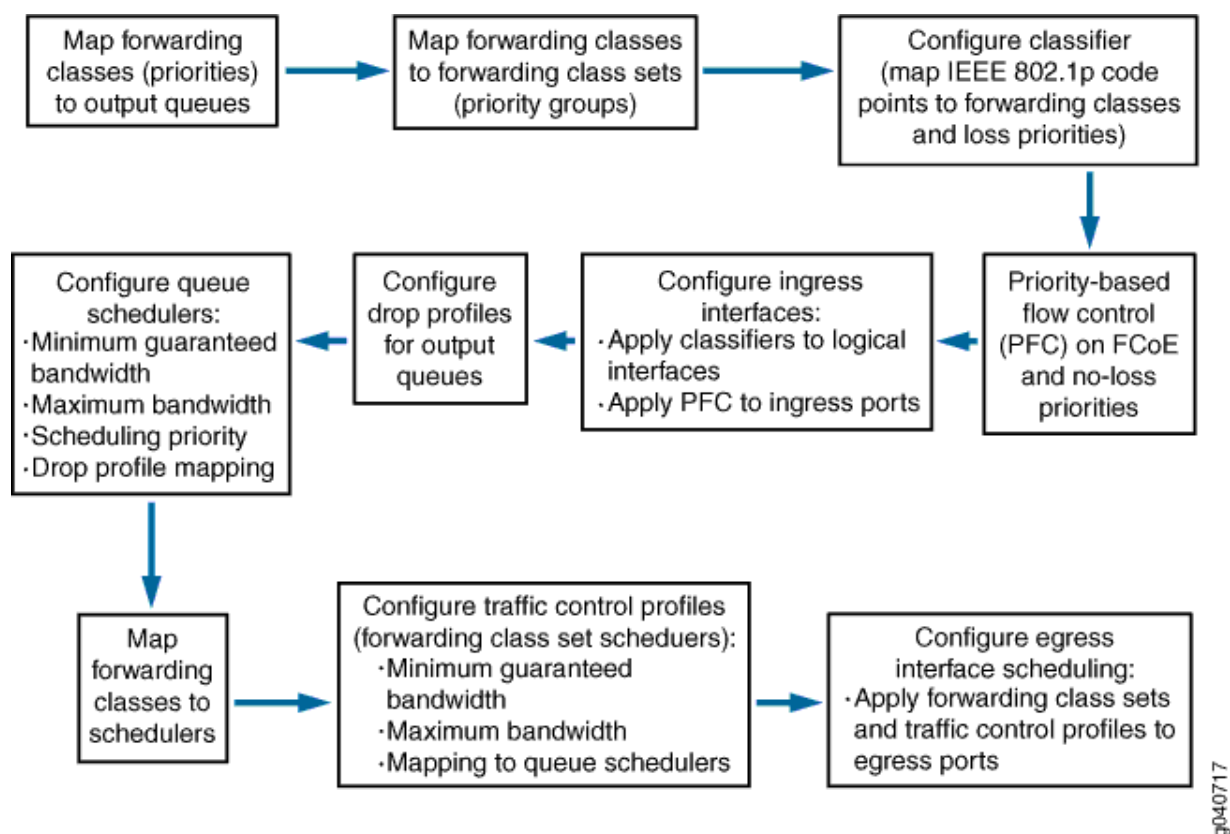
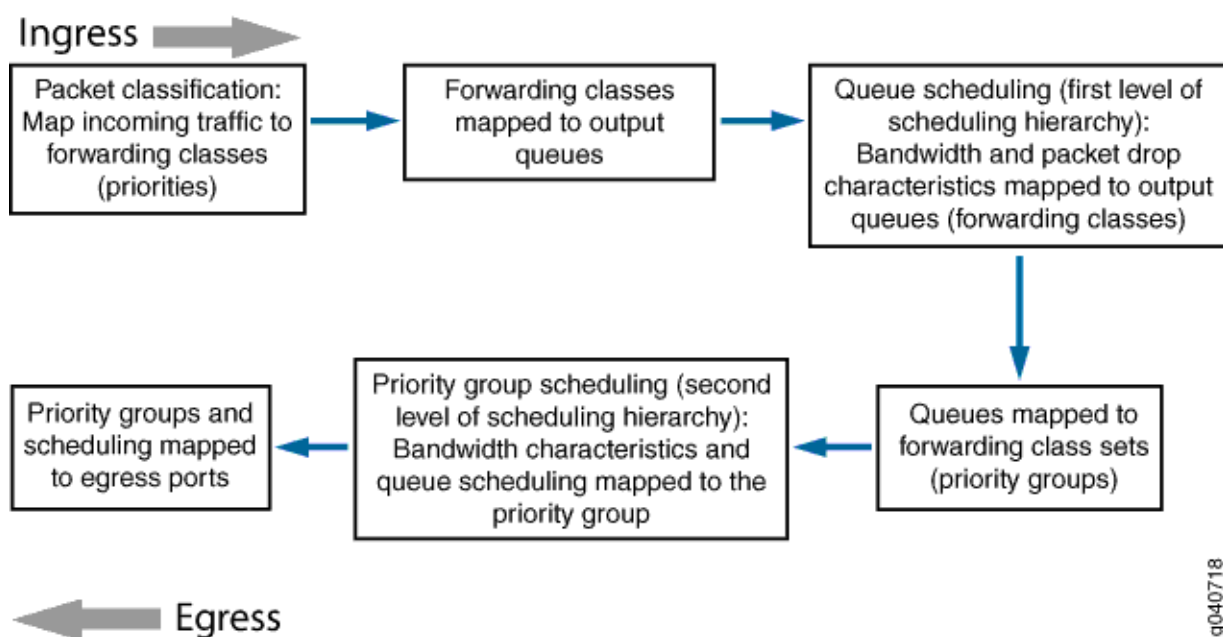


Figure 22 on page 452 shows a block diagram of the hierarchical scheduling packet flow from ingress to egress.

Figure 22: Hierarchical Port Scheduling Packet Flow Block Diagram



Configuration

IN THIS SECTION

- [CLI Quick Configuration | 452](#)
- [Procedure | 456](#)
- [Results | 462](#)

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 \(Priorities\) | 467](#)
- [Verifying the Forwarding Class Sets \(Priority Groups\) | 468](#)
- [Verifying the Classifier | 469](#)
- [Verifying Priority-Based Flow Control | 470](#)

- [Verifying the Output Queue Schedulers | 471](#)
- [Verifying the Drop Profiles | 475](#)
- [Verifying the Priority Group Output Schedulers \(Traffic Control Profiles\) | 476](#)
- [Verifying the Interface Configuration | 477](#)

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: 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 \(DSCP, DSCP IPv6, IEEE 802.1p\) | 106](#)

[Benefits of Configuring CoS Hierarchical Port Scheduling](#)

[Assigning CoS Components to Interfaces | 87](#)

[Example: Configuring WRED Drop Profiles | 286](#)

[Example: Configuring Drop Profile Maps | 293](#)

[Example: Configuring Forwarding Classes | 174](#)

[Example: Configuring Forwarding Class Sets | 184](#)

[Example: Configuring Queue Schedulers | 350](#)

[Example: Configuring Queue Scheduling Priority | 360](#)

[Example: Configuring Traffic Control Profiles \(Priority Group Scheduling\) | 414](#)

[Example: Configuring Minimum Guaranteed Output Bandwidth | 421](#)

[Example: Configuring Maximum Output Bandwidth | 431](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

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

[Understanding CoS Hierarchical Port Scheduling \(ETS\) | 438](#)

[Understanding CoS Scheduling Behavior and Configuration Considerations | 332](#)

[*Understanding CoS Scheduling on QFabric System Node Device Fabric \(fte\) Ports*](#)

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

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 | 496](#)

[Configuring DCBX Autonegotiation | 497](#)

[Understanding DCBX | 486](#)

[Understanding Data Center Bridging Capability Exchange Protocol for EX Series Switches](#)

# 4

PART

## Data Center Bridging and Lossless FCoE

---

Data Center Bridging | 482

Lossless FCoE | 524

---

# Data Center Bridging

## IN THIS CHAPTER

- [Understanding DCB Features and Requirements | 482](#)
- [Understanding DCBX | 486](#)
- [Configuring the DCBX Mode | 496](#)
- [Configuring DCBX Autonegotiation | 497](#)
- [Understanding DCBX Application Protocol TLV Exchange | 500](#)
- [Defining an Application for DCBX Application Protocol TLV Exchange | 504](#)
- [Configuring an Application Map for DCBX Application Protocol TLV Exchange | 506](#)
- [Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange | 507](#)
- [Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

## Understanding DCB Features and Requirements

### IN THIS SECTION

- [Lossless Transport | 483](#)
- [ETS | 484](#)
- [DCBX | 485](#)

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 \(Ethernet PAUSE and PFC\) | 220](#)

[Understanding DCBX | 486](#)

[Example: Configuring CoS PFC for FCoE Traffic | 524](#)

## Understanding DCBX

### IN THIS SECTION

- [DCBX Basics | 486](#)
- [DCBX Modes and Support | 487](#)
- [DCBX Attribute Types | 490](#)
- [DCBX Application Protocol TLV Exchange | 492](#)
- [DCBX and PFC | 493](#)
- [DCBX and ETS | 493](#)

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 85 on page 488 summarizes the differences between IEEE DCBX and DCBX version 1.01, including show command output:

**Table 85: Summary of Differences Between IEEE DCBX and DCBX Version 1.01**

| Characteristic                               | IEEE DCBX                                                                                                                                                                                                  | DCBX Version 1.01                                                                                                                                                                                    |
|----------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| OUI                                          | 0x0080c2                                                                                                                                                                                                   | 0x001b21                                                                                                                                                                                             |
| 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                                                                                                                                                                                       |

**Table 85: Summary of Differences Between IEEE DCBX and DCBX Version 1.01 (Continued)**

| Characteristic                                                                                | IEEE DCBX                                                                                                                                                                                                                                                                                                                                                                                                                                                                        | DCBX Version 1.01                                                                                                                                                                                                                                                                                                                                                                                                                                                                                            |
|-----------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| Differences in the <code>show dcbx interface <i>interface-name</i></code> operational command | <ul style="list-style-type: none"> <li>• Synchronization information is not shown because symmetric configuration is not required.</li> <li>• Operational state information is not shown because the operational states do not have to be symmetric.</li> <li>• TLV type is shown because unique TLVs are sent for each DCBX attribute.</li> <li>• ETS peer Configuration TLV and Recommendation TLV information is shown separately because they are different TLVs.</li> </ul> | <ul style="list-style-type: none"> <li>• Synchronization information is shown because symmetric configuration is required.</li> <li>• Operational state information is shown because the operational states do have to be symmetric.</li> <li>• TLV type is not shown because one TLV is used for all attribute information.</li> <li>• 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).</li> </ul> |

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 | 500](#)

[Understanding DCB Features and Requirements | 482](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

[Understanding CoS Hierarchical Port Scheduling \(ETS\)](#)

[Understanding CoS Port Schedulers on QFX Switches](#)

[Understanding FCoE](#)

[Configuring the DCBX Mode | 496](#)

[Configuring DCBX Autonegotiation | 497](#)

[Disabling the ETS Recommendation TLV | 480](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

## 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 | 497](#)

[Disabling the ETS Recommendation TLV | 480](#)

[Understanding DCBX | 486](#)

[Understanding DCBX Application Protocol TLV Exchange | 500](#)

*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 | 509](#)

[Example: Configuring CoS PFC for FCoE Traffic | 524](#)

[Disabling the ETS Recommendation TLV | 480](#)

[Understanding DCBX Application Protocol TLV Exchange | 500](#)

Understanding DCBX Application Protocol TLV Exchange

IN THIS SECTION

- [Applications | 501](#)
- [Application Maps | 502](#)
- [Classifying and Prioritizing Application Traffic | 503](#)
- [Enabling Interfaces to Exchange Application Protocol Information | 503](#)
- [Disabling DCBX Application Protocol Exchange | 504](#)

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 501](#)) 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 502](#)).

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 | 486](#)

[Configuring DCBX Autonegotiation | 497](#)

[Disabling the ETS Recommendation TLV | 480](#)

[Defining an Application for DCBX Application Protocol TLV Exchange | 504](#)

[Configuring an Application Map for DCBX Application Protocol TLV Exchange | 506](#)

[Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange | 507](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

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 | 506](#)

[Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange | 507](#)

[Configuring DCBX Autonegotiation | 497](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

[Example: Configuring DCBX to Support an iSCSI Application](#)

[Understanding DCBX Application Protocol TLV Exchange | 500](#)

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 | 504](#)

[Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange | 507](#)

[Configuring DCBX Autonegotiation | 497](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

[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 | 504](#)

[Configuring an Application Map for DCBX Application Protocol TLV Exchange | 506](#)

[Configuring DCBX Autonegotiation | 497](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

[Example: Configuring DCBX to Support an iSCSI Application](#)

show dcbx neighbors

Example: Configuring DCBX Application Protocol TLV Exchange

IN THIS SECTION

- [Requirements | 510](#)
- [Overview | 510](#)
- [Configuration | 515](#)
- [Verification | 518](#)

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 | 512](#)

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" on page 497](#) 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 86 on page 511](#) 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 87 on page 511](#) shows the default untrusted classifier IEEE 802.1 code-point values to unicast forwarding class mapping for ports in access mode.

Table 86: 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 87: 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 87: 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 88 on page 512](#) shows the configuration components for this example.

Table 88: Components of DCBX Application Protocol Exchange Configuration Topology

Component	Settings
Hardware	QFX Series device

Table 88: 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 protocol—TCP destination-port—3260 code-points—111
PTP application (Layer 2)	Application name—ptp ether-type—0x88F7 code-points—001, 101
FCoE application (Layer 2)	Application name—fcoe ether-type—0x8906 code-points—011 NOTE: 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 dcbx-iscsi-ptp-app-map—Maps iSCSI and PTP applications to IEEE 802.1p code points

Table 88: 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-ntp-app-map)</p>
PFC congestion notification profile for FCoE application exchange	<p>fcoe-cnp:</p> <ul style="list-style-type: none"> • Code point—011 • Interface—xe-0/0/10
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"> • Maps the fcoe forwarding class to the IEEE 802.1p code point used for the FCoE application (011) and a loss priority of high • 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 • Applied to interface xe-0/0/10 <p>iscsi-ntp-cl2:</p> <ul style="list-style-type: none"> • 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 • 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 • Applied to interface xe-0/0/11

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 | 515](#)
- [Configuring DCBX Application Protocol TLV Exchange | 516](#)

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 | 518](#)
- [Verifying the Application Map Configuration | 519](#)
- [Verifying DCBX Application Protocol Exchange Interface Configuration | 520](#)
- [Verifying the PFC Configuration | 520](#)
- [Verifying the Classifier Configuration | 521](#)

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-ntp-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-ntp-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-ntp-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 | 504](#)

[Configuring an Application Map for DCBX Application Protocol TLV Exchange | 506](#)

[Applying an Application Map to an Interface for DCBX Application Protocol TLV Exchange | 507](#)

[Configuring DCBX Autonegotiation | 497](#)

show dcbx

show dcbx neighbors

[Understanding DCBX Application Protocol TLV Exchange | 500](#)

Lossless FCoE

IN THIS CHAPTER

- [Example: Configuring CoS PFC for FCoE Traffic | 524](#)
- [Example: Configuring CoS for FCoE Transit Switch Traffic Across an MC-LAG | 538](#)
- [Example: Configuring CoS Using ELS for FCoE Transit Switch Traffic Across an MC-LAG | 572](#)
- [Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic \(FCoE Transit Switch\) | 608](#)
- [Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface | 620](#)
- [Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces | 633](#)
- [Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications \(FCoE and iSCSI\) | 652](#)
- [Troubleshooting Dropped FCoE Traffic | 678](#)

Example: Configuring CoS PFC for FCoE Traffic

IN THIS SECTION

- [Requirements | 525](#)
- [Overview | 525](#)
- [Configuration | 528](#)
- [Verification | 535](#)

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 | 526](#)

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 89 on page 526](#) shows the configuration components for this example.

Table 89: Components of the PFC for FCoE Traffic Configuration Topology

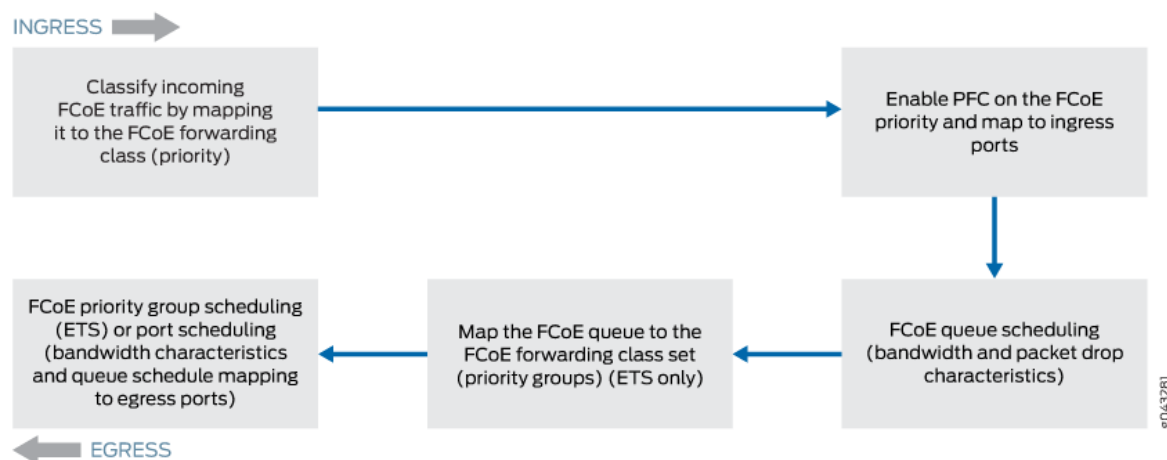
Component	Settings
Hardware	One switch
Behavior aggregate classifier (maps the FCoE forwarding class to incoming packets by IEEE 802.1 code point)	Code point 011 to forwarding class fcoe and loss priority low Ingress interfaces: xe-0/0/31, xe-0/0/32, xe-0/0/33, xe-0/0/34
PFC congestion notification profile	fcoe-cnp: Code point 011 Ingress interfaces: xe-0/0/31, xe-0/0/32, xe-0/0/33, xe-0/0/34

Table 89: Components of the PFC for FCoE Traffic Configuration Topology (*Continued*)

Component	Settings
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 23 on page 528 shows a block diagram of the configuration components and the configuration flow of the CLI statements used in the example.

Figure 23: PFC for FCoE Traffic Configuration Components Block Diagram



Configuration

IN THIS SECTION

- [CLI Quick Configuration | 528](#)
- [Common Configuration \(Applies to ETS Hierarchical Scheduling and to Port Scheduling\) | 530](#)
- [ETS Hierarchical Scheduling Configuration | 531](#)
- [Port Scheduling Configuration | 532](#)
- [Results | 532](#)

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 | 535](#)
- [Verifying the Ingress Interface PFC Configuration | 536](#)

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 \(Ethernet PAUSE and PFC\) | 220](#)

Example: Configuring CoS for FCoE Transit Switch Traffic Across an MC-LAG

IN THIS SECTION

- [Requirements | 539](#)
- [Overview | 539](#)
- [Configuration | 546](#)
- [Verification | 559](#)

Multichassis link aggregation groups (MC-LAGs) provide redundancy and load balancing between two switches, multihoming support for client devices such as servers, and a loop-free Layer 2 network without running Spanning Tree Protocol (STP).

NOTE: This example uses Junos OS without support for the Enhanced Layer 2 Software (ELS) configuration style. If your switch runs software that does support ELS, see ["Example: Configuring CoS Using ELS for FCoE Transit Switch Traffic Across an MC-LAG" on page 572](#). For ELS details, see [Using the Enhanced Layer 2 Software CLI](#).

You can use an MC-LAG to provide a redundant aggregation layer for Fibre Channel over Ethernet (FCoE) traffic in an *inverted-U* topology. To support lossless transport of FCoE traffic across an MC-LAG, you must configure the appropriate class of service (CoS) on both of the switches with MC-LAG port

members. The CoS configuration must be the same on both of the MC-LAG switches because an MC-LAG does not carry forwarding class and IEEE 802.1p priority information.

NOTE: This example describes how to configure CoS to provide lossless transport for FCoE traffic across an MC-LAG that connects two switches. It also describes how to configure CoS on the FCoE transit switches that connect FCoE hosts to the two switches that form the MC-LAG. This example does *not* describe how to configure the MC-LAG itself. However, this example includes a subset of MC-LAG configuration that only shows how to configure interface membership in the MC-LAG.

Ports that are part of an FCoE-FC gateway configuration (a virtual FCoE-FC gateway fabric) do not support MC-LAGs. Ports that are members of an MC-LAG act as FCoE pass-through transit switch ports.

QFX Series switches and EX4600 switches support MC-LAGs. QFabric system Node devices do not support MC-LAGs.

Requirements

This example uses the following hardware and software components:

- Two Juniper Networks QFX3500 switches that form an MC-LAG for FCoE traffic.
- Two Juniper Networks QFX3500 switches that provide FCoE server access in transit switch mode and that connect to the MC-LAG switches. These switches can be standalone QFX3500 switches or they can be Node devices in a QFabric system.
- FCoE servers (or other FCoE hosts) connected to the transit switches.
- Junos OS Release 12.2 or later for the QFX Series.

Overview

IN THIS SECTION

- [Topology | 540](#)

FCoE traffic requires lossless transport. This example shows you how to:

- Configure CoS for FCoE traffic on the two QFX3500 switches that form the MC-LAG, including priority-based flow control (PFC) and enhanced transmission selection (ETS; hierarchical scheduling of resources for the FCoE forwarding class priority and for the forwarding class set priority group).

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 CoS for FCoE on the two FCoE transit switches that connect FCoE hosts to the MC-LAG switches and enable FIP snooping on the FCoE VLAN at the FCoE transit switch access ports.
- Disable IGMP snooping on the FCoE VLAN.

NOTE: This is only necessary if IGMP snooping is enabled on the VLAN. Before Junos OS Release 13.2, IGMP snooping was enabled by default on VLANs. Beginning with Junos OS Release 13.2, IGMP snooping is enabled by default only on the default VLAN.

- Configure the appropriate port mode, MTU, and FCoE trusted or untrusted state for each interface to support lossless FCoE transport.

Topology

Switches that act as transit switches support MC-LAGs for FCoE traffic in an inverted-U network topology, as shown in [Figure 24 on page 541](#).

Figure 24: Supported Topology for an MC-LAG on an FCoE Transit Switch

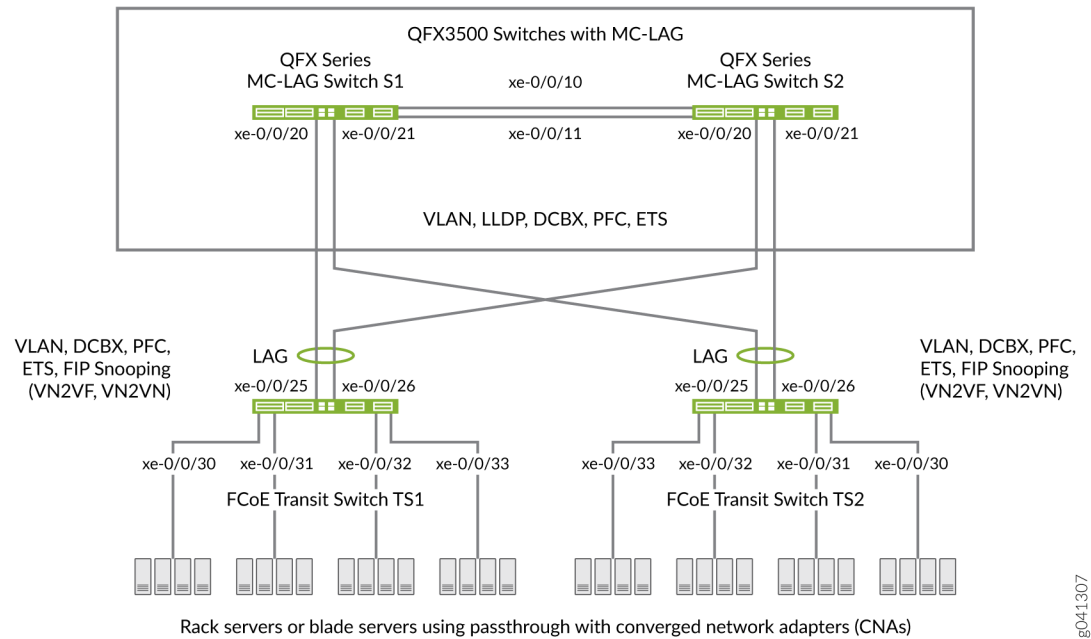


Table 90 on page 541 shows the configuration components for this example.

Table 90: Components of the CoS for FCoE Traffic Across an MC-LAG Configuration Topology

Component	Settings
Hardware	Four QFX3500 switches (two to form the MC-LAG as pass-through transit switches and two transit switches for FCoE access).
Forwarding class (all switches)	Default fcoe forwarding class.
Classifier (forwarding class mapping of incoming traffic to IEEE priority)	Default IEEE 802.1p trusted classifier on all FCoE interfaces.

Table 90: Components of the CoS for FCoE Traffic Across an MC-LAG Configuration Topology
(Continued)

Component	Settings
LAGs and MC-LAG	<p>S1—Ports xe-0/0/10 and x-0/0/11 are members of LAG ae0, which connects Switch S1 to Switch S2. Ports xe-0/0/20 and xe-0/0/21 are members of MC-LAG ae1. All ports are configured in trunk port mode, as fcoe-trusted, and with an MTU of 2180.</p> <p>S2—Ports xe-0/0/10 and x-0/0/11 are members of LAG ae0, which connects Switch S2 to Switch S1. Ports xe-0/0/20 and xe-0/0/21 are members of MC-LAG ae1. All ports are configured in trunk port mode, as fcoe-trusted, and with an MTU of 2180.</p> <p>NOTE: Ports xe-0/0/20 and xe-0/0/21 on Switches S1 and S2 are the members of the MC-LAG.</p> <p>TS1—Ports xe-0/0/25 and x-0/0/26 are members of LAG ae1, configured in trunk port mode, as fcoe-trusted, and with an MTU of 2180. Ports xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 are configured in tagged-access port mode, with an MTU of 2180.</p> <p>TS2—Ports xe-0/0/25 and x-0/0/26 are members of LAG ae1, configured in trunk port mode, as fcoe-trusted, and with an MTU of 2180. Ports xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 are configured in tagged-access port mode, with an MTU of 2180.</p>
FCoE queue scheduler (all switches)	<p>fcoe-sched: Minimum bandwidth 3g Maximum bandwidth 100% Priority low</p>

Table 90: Components of the CoS for FCoE Traffic Across an MC-LAG Configuration Topology
(Continued)

Component	Settings
Forwarding class-to-scheduler mapping (all switches)	Scheduler map fcoe-map: Forwarding class fcoe Scheduler fcoe-sched
Forwarding class set (FCoE priority group, all switches)	fcoe-pg: Forwarding class fcoe Egress interfaces: <ul style="list-style-type: none"> • S1—LAG ae0 and MC-LAG ae1 • S2—LAG ae0 and MC-LAG ae1 • TS1—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 • TS2—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33
Traffic control profile (all switches)	fcoe-tcp: Scheduler map fcoe-map Minimum bandwidth 3g Maximum bandwidth 100%
PFC congestion notification profile (all switches)	fcoe-cnp: Code point 011 Ingress interfaces: <ul style="list-style-type: none"> • S1—LAG ae0 and MC-LAG ae1 • S2—LAG ae0 and MC-LAG ae1 • TS1—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 • TS2—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33

Table 90: Components of the CoS for FCoE Traffic Across an MC-LAG Configuration Topology
(Continued)

Component	Settings
FCoE VLAN name and tag ID	<p>Name—fcoe_vlan ID—100</p> <p>Include the FCoE VLAN on the interfaces that carry FCoE traffic on all four switches.</p> <p>Disable IGMP snooping on the interfaces that belong to the FCoE VLAN on all four switches.</p>
FIP snooping	<p>Enable FIP snooping on Transit Switches TS1 and TS2 on the FCoE VLAN. Configure the LAG interfaces that connect to the MC-LAG switches as FCoE trusted interfaces so that they do not perform FIP snooping.</p> <p>This example enables VN2VN_Port FIP snooping on the FCoE transit switch interfaces connected to the FCoE servers. The example is equally valid with VN2VF_Port FIP snooping enabled on the transit switch access ports. The method of FIP snooping you enable depends on your network configuration.</p>

NOTE: This example uses the default IEEE 802.1p trusted BA classifier, which is automatically applied to trunk mode and tagged access mode ports if you do not apply an explicitly configured classifier.

To configure CoS for FCoE traffic across an MC-LAG:

- Use the default FCoE forwarding class and forwarding-class-to-queue mapping (do not explicitly configure the FCoE forwarding class or output queue). The default FCoE forwarding class is `fcoe`, and the default output queue is `queue 3`.

NOTE: In Junos OS Release 12.2, traffic mapped to explicitly configured forwarding classes, even lossless forwarding classes such as `fcoe`, is treated as lossy (best-effort) traffic and does *not* receive lossless treatment. To receive lossless treatment in Release 12.2, traffic must use one of the default lossless forwarding classes (`fcoe` or `no-loss`).

In Junos OS Release 12.3 and later, you can include the *no-loss* packet drop attribute in the explicit forwarding class configuration to configure a lossless forwarding class.

- Use the default trusted BA classifier, which maps incoming packets to forwarding classes by the IEEE 802.1p code point (CoS priority) of the packet. The trusted classifier is the default classifier for interfaces in trunk and tagged-access port modes. The default trusted classifier maps incoming packets with the IEEE 802.1p code point 3 (011) to the FCoE forwarding class. If you choose to configure the BA classifier instead of using the default classifier, you must ensure that FCoE traffic is classified into forwarding classes in exactly the same way on both MC-LAG switches. Using the default classifier ensures consistent classifier configuration on the MC-LAG ports.
- Configure a congestion notification profile that enables PFC on the FCoE code point (code point 011 in this example). The congestion notification profile configuration must be the same on both MC-LAG switches.
- Apply the congestion notification profile to the interfaces.
- Configure enhanced transmission selection (ETS, also known as hierarchical scheduling) on the interfaces to provide the bandwidth required for lossless FCoE transport. Configuring ETS includes configuring bandwidth scheduling for the FCoE forwarding class, a forwarding class set (priority group) that includes the FCoE forwarding class, and a traffic control profile to assign bandwidth to the forwarding class set that includes FCoE traffic.
- Apply the ETS scheduling to the interfaces.
- Configure the port mode, MTU, and FCoE trusted or untrusted state for each interface to support lossless FCoE transport.

In addition, this example describes how to enable FIP snooping on the Transit Switch TS1 and TS2 ports that are connected to the FCoE servers and how to disable IGMP snooping on the FCoE VLAN. To provide secure access, FIP snooping must be enabled on the FCoE access ports.

This example focuses on the CoS configuration to support lossless FCoE transport across an MC-LAG. This example does not describe how to configure the properties of MC-LAGs and LAGs, although it does show you how to configure the port characteristics required to support lossless transport and how to assign interfaces to the MC-LAG and to the LAGs.

Before you configure CoS, configure:

- The MC-LAGs that connect Switches S1 and S2 to Switches TS1 and TS2.
- The LAGs that connect the Transit Switches TS1 and TS2 to MC-LAG Switches S1 and S2.
- The LAG that connects Switch S1 to Switch S2.

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 546](#)
- [Configuring MC-LAG Switches S1 and S2 | 548](#)
- [Configuring FCoE Transit Switches TS1 and TS2 | 551](#)
- [Results | 555](#)

To configure CoS for lossless FCoE transport across an MC-LAG, perform these tasks:

CLI Quick Configuration

To quickly configure CoS for lossless FCoE transport across an MC-LAG, 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 MC-LAG Switch S1 and MC-LAG Switch S2 at the [edit] hierarchy level. The configurations on Switches S1 and S2 are identical because the CoS configuration must be identical, and because this example uses the same ports on both switches.

Switch S1 and Switch S2

```
set class-of-service schedulers fcoe-sched priority low transmit-rate 3g
set class-of-service schedulers fcoe-sched shaping-rate percent 100
set class-of-service scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
set class-of-service forwarding-class-sets fcoe-pg class fcoe
set class-of-service traffic-control-profiles fcoe-tcp scheduler-map fcoe-map guaranteed-rate 3g
set class-of-service traffic-control-profiles fcoe-tcp shaping-rate percent 100
set class-of-service interfaces ae0 forwarding-class-set fcoe-pg output-traffic-control-profile
fcoe-tcp
set class-of-service interfaces ae1 forwarding-class-set fcoe-pg output-traffic-control-profile
fcoe-tcp
set class-of-service congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
set class-of-service interfaces ae0 congestion-notification-profile fcoe-cnp
set class-of-service interfaces ae1 congestion-notification-profile fcoe-cnp
set vlans fcoe_vlan vlan-id 100
set protocols igmp-snooping vlan fcoe_vlan disable
set interfaces xe-0/0/10 ether-options 802.3ad ae0
```

```

set interfaces xe-0/0/11 ether-options 802.3ad ae0
set interfaces xe-0/0/20 ether-options 802.3ad ae1
set interfaces xe-0/0/21 ether-options 802.3ad ae1
set interfaces ae0 unit 0 family ethernet-switching port-mode trunk vlan members fcoe_vlan
set interfaces ae1 unit 0 family ethernet-switching port-mode trunk vlan members fcoe_vlan
set interfaces ae0 mtu 2180
set interfaces ae1 mtu 2180
set ethernet-switching-options secure-access-port interface ae0 fcoe-trusted
set ethernet-switching-options secure-access-port interface ae1 fcoe-trusted

```

To quickly configure CoS for lossless FCoE transport across an MC-LAG, 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 Transit Switch TS1 and Transit Switch TS2 at the [edit] hierarchy level. The configurations on Switches TS1 and TS2 are identical because the CoS configuration must be identical, and because this example uses the same ports on both switches.

Switch TS1 and Switch TS2

```

set class-of-service schedulers fcoe-sched priority low transmit-rate 3g
set class-of-service schedulers fcoe-sched shaping-rate percent 100
set class-of-service scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
set class-of-service forwarding-class-sets fcoe-pg class fcoe
set class-of-service traffic-control-profiles fcoe-tcp scheduler-map fcoe-map guaranteed-rate 3g
set class-of-service traffic-control-profiles fcoe-tcp shaping-rate percent 100
set class-of-service interfaces ae1 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service interfaces xe-0/0/30 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service interfaces xe-0/0/31 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service interfaces xe-0/0/32 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service interfaces xe-0/0/33 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
set class-of-service interfaces ae1 congestion-notification-profile fcoe-cnp
set class-of-service interfaces xe-0/0/30 congestion-notification-profile fcoe-cnp
set class-of-service interfaces xe-0/0/31 congestion-notification-profile fcoe-cnp
set class-of-service interfaces xe-0/0/32 congestion-notification-profile fcoe-cnp
set class-of-service interfaces xe-0/0/33 congestion-notification-profile fcoe-cnp
set vlans fcoe_vlan vlan-id 100

```



```

set protocols igmp-snooping vlan fcoe_vlan disable
set interfaces xe-0/0/25 ether-options 802.3ad ae1
set interfaces xe-0/0/26 ether-options 802.3ad ae1
set interfaces ae1 unit 0 family ethernet-switching port-mode trunk vlan members fcoe_vlan
set interfaces xe-0/0/30 unit 0 family ethernet-switching port-mode tagged-access vlan members fcoe_vlan
set interfaces xe-0/0/31 unit 0 family ethernet-switching port-mode tagged-access vlan members fcoe_vlan
set interfaces xe-0/0/32 unit 0 family ethernet-switching port-mode tagged-access vlan members fcoe_vlan
set interfaces xe-0/0/33 unit 0 family ethernet-switching port-mode tagged-access vlan members fcoe_vlan
set interfaces ae1 mtu 2180
set interfaces xe-0/0/30 mtu 2180
set interfaces xe-0/0/31 mtu 2180
set interfaces xe-0/0/32 mtu 2180
set interfaces xe-0/0/33 mtu 2180
set ethernet-switching-options secure-access-port interface ae1 fcoe-trusted
set ethernet-switching-options secure-access-port vlan fcoe_vlan examine-fip examine-vn2v2
beacon-period 90000

```

Configuring MC-LAG Switches S1 and S2

Step-by-Step Procedure

To configure CoS resource scheduling (ETS), PFC, the FCoE VLAN, and the LAG and MC-LAG interface membership and characteristics to support lossless FCoE transport across an MC-LAG (this example uses the default `fcoe` forwarding class and the default classifier to map incoming FCoE traffic to the FCoE IEEE 802.1p code point 011, so you do not configure them):

1. Configure output scheduling for the FCoE queue.

```

[edit class-of-service schedulers fcoe-sched]
user@switch# set priority low transmit-rate 3g
user@switch# set shaping-rate percent 100

```

2. Map the FCoE forwarding class to the FCoE scheduler (fcoe-sched).

```
[edit class-of-service]
user@switch# set scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
```

3. Configure the forwarding class set (fcoe-pg) for the FCoE traffic.

```
[edit class-of-service]
user@switch# set forwarding-class-sets fcoe-pg class fcoe
```

4. Define the traffic control profile (fcoe-tcp) to use on the FCoE forwarding class set.

```
[edit class-of-service traffic-control-profiles fcoe-tcp]
user@switch# set scheduler-map fcoe-map guaranteed-rate 3g
user@switch# set shaping-rate percent 100
```

5. Apply the FCoE forwarding class set and traffic control profile to the LAG and MC-LAG interfaces.

```
[edit class-of-service]
user@switch# set interfaces ae0 forwarding-class-set fcoe-pg output-traffic-control-profile
fcoe-tcp
user@switch# set interfaces ae1 forwarding-class-set fcoe-pg output-traffic-control-profile
fcoe-tcp
```

6. Enable PFC on the FCoE priority by creating a congestion notification profile (fcoe-cnp) that applies 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
```

7. Apply the PFC configuration to the LAG and MC-LAG interfaces.

```
[edit class-of-service]
user@switch# set interfaces ae0 congestion-notification-profile fcoe-cnp
user@switch# set interfaces ae1 congestion-notification-profile fcoe-cnp
```

8. Configure the VLAN for FCoE traffic (fcoe_vlan).

```
[edit vlans]
user@switch# set fcoe_vlan vlan-id 100
```

9. Disable IGMP snooping on the FCoE VLAN.

```
[edit protocols]
user@switch# set igmp-snooping vlan fcoe_vlan disable
```

10. Add the member interfaces to the LAG between the two MC-LAG switches.

```
[edit interfaces]
user@switch# set xe-0/0/10 ether-options 802.3ad ae0
user@switch# set xe-0/0/11 ether-options 802.3ad ae0
```

11. Add the member interfaces to the MC-LAG.

```
[edit interfaces]
user@switch# set xe-0/0/20 ether-options 802.3ad ae1
user@switch# set xe-0/0/21 ether-options 802.3ad ae1
```

12. Configure the port mode as trunk and membership in the FCoE VLAN (fcoe_vlan) for the LAG (ae0) and for the MC-LAG (ae1).

```
[edit interfaces]
user@switch# set ae0 unit 0 family ethernet-switching port-mode trunk vlan members fcoe_vlan
user@switch# set ae1 unit 0 family ethernet-switching port-mode trunk vlan members fcoe_vlan
```

13. Set the MTU to 2180 for the LAG and MC-LAG interfaces.

2180 bytes is the minimum size required to handle FCoE packets because of the payload and header sizes. You can configure the MTU to a higher number of bytes if desired, but not less than 2180 bytes.

```
[edit interfaces]
user@switch# set ae0 mtu 2180
user@switch# set ae1 mtu 2180
```

14. Set the LAG and MC-LAG interfaces as FCoE trusted ports.

Ports that connect to other switches should be trusted and should not perform FIP snooping.

```
[edit ethernet-switching-options secure-access-port interface]
user@switch# set ae0 fcoe-trusted
user@switch# set ae1 fcoe-trusted
```

Configuring FCoE Transit Switches TS1 and TS2

Step-by-Step Procedure

The CoS configuration on FCoE Transit Switches TS1 and TS2 is similar to the CoS configuration on MC-LAG Switches S1 and S2. However, the port configurations differ, and you must enable FIP snooping on the Switch TS1 and Switch TS2 FCoE access ports.

To configure resource scheduling (ETS), PFC, the FCoE VLAN, and the LAG interface membership and characteristics to support lossless FCoE transport across the MC-LAG (this example uses the default `fcoe` forwarding class and the default classifier to map incoming FCoE traffic to the FCoE IEEE 802.1p code point 011, so you do not configure them):

1. Configure output scheduling for the FCoE queue.

```
[edit class-of-service schedulers fcoe-sched]
user@switch# set priority low transmit-rate 3g
user@switch# set shaping-rate percent 100
```

2. Map the FCoE forwarding class to the FCoE scheduler (fcoe-sched).

```
[edit class-of-service]
user@switch# set scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
```

3. Configure the forwarding class set (fcoe-pg) for the FCoE traffic.

```
[edit class-of-service]
user@switch# set forwarding-class-sets fcoe-pg class fcoe
```

4. Define the traffic control profile (fcoe-tcp) to use on 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
```

5. Apply the FCoE forwarding class set and traffic control profile to the LAG interface and to the FCoE access interfaces.

```
[edit class-of-service]
user@switch# set interfaces ae1 forwarding-class-set fcoe-pg output-traffic-control-profile
fcoe-tcp
user@switch# set interfaces xe-0/0/30 forwarding-class-set fcoe-pg output-traffic-control-
profile fcoe-tcp
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
```

6. Enable PFC on the FCoE priority by creating a congestion notification profile (fcoe-cnp) that applies 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
```

7. Apply the PFC configuration to the LAG interface and to the FCoE access interfaces.

```
[edit class-of-service]
user@switch# set interfaces ae1 congestion-notification-profile fcoe-cnp
user@switch# set interfaces xe-0/0/30 congestion-notification-profile fcoe-cnp
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
```

8. Configure the VLAN for FCoE traffic (fcoe_vlan).

```
[edit vlans]
user@switch# set fcoe_vlan vlan-id 100
```

9. Disable IGMP snooping on the FCoE VLAN.

```
[edit protocols]
user@switch# set igmp-snooping vlan fcoe_vlan disable
```

10. Add the member interfaces to the LAG.

```
[edit interfaces]
user@switch# set xe-0/0/25 ether-options 802.3ad ae1
user@switch# set xe-0/0/26 ether-options 802.3ad ae1
```

11. On the LAG (ae1), configure the port mode as trunk and membership in the FCoE VLAN (fcoe_vlan).

```
[edit interfaces]
user@switch# set ae1 unit 0 family ethernet-switching port-mode trunk vlan members fcoe_vlan
```

12. On the FCoE access interfaces (xe-0/0/30, xe-0/0/31, xe-0/0/32, xe-0/0/33), configure the port mode as tagged-access and membership in the FCoE VLAN (fcoe_vlan).

```
[edit interfaces]
user@switch# set xe-0/0/30 unit 0 family ethernet-switching port-mode tagged-access vlan
members fcoe_vlan
user@switch# set xe-0/0/31 unit 0 family ethernet-switching port-mode tagged-access vlan
members fcoe_vlan
user@switch# set xe-0/0/32 unit 0 family ethernet-switching port-mode tagged-access vlan
members fcoe_vlan
user@switch# set xe-0/0/33 unit 0 family ethernet-switching port-mode tagged-access vlan
members fcoe_vlan
```

13. Set the MTU to 2180 for the LAG and FCoE access interfaces.

2180 bytes is the minimum size required to handle FCoE packets because of the payload and header sizes; you can configure the MTU to a higher number of bytes if desired, but not less than 2180 bytes.

```
[edit interfaces]
user@switch# set ae1 mtu 2180
user@switch# set xe-0/0/30 mtu 2180
user@switch# set xe-0/0/31 mtu 2180
user@switch# set xe-0/0/32 mtu 2180
user@switch# set xe-0/0/33 mtu 2180
```

14. Set the LAG interface as an FCoE trusted port. Ports that connect to other switches should be trusted and should not perform FIP snooping:

```
[edit ethernet-switching-options]
user@switch# set secure-access-port interface ae1 fcoe-trusted
```

NOTE: Access ports xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 are not configured as FCoE trusted ports. The access ports remain in the default state as untrusted ports because they connect directly to FCoE devices and must perform FIP snooping to ensure network security.

15. Enable FIP snooping on the FCoE VLAN to prevent unauthorized FCoE network access (this example uses VN2VN_Port FIP snooping; the example is equally valid if you use VN2VF_Port FIP snooping).

```
[edit ethernet-switching-options]
user@switch# set secure-access-port vlan fcoe_vlan examine-fip examine-vn2vn beacon-period
90000
```

Results

Display the results of the CoS configuration on MC-LAG Switch S1 and on MC-LAG Switch S2 (the results on both switches are the same).

```
user@switch> show configuration class-of-service
traffic-control-profiles {
  fcoe-tcp {
    scheduler-map fcoe-map;
    shaping-rate percent 100;
    guaranteed-rate 3g;
  }
}
forwarding-class-sets {
  fcoe-pg {
    class fcoe;
  }
}
congestion-notification-profile {
  fcoe-cnp {
    input {
      ieee-802.1 {
        code-point 011 {
          pfc;
        }
      }
    }
  }
}
```



```

    }
  }
}
interfaces {
  ae0 {
    forwarding-class-set {
      fcoe-pg {
        output-traffic-control-profile fcoe-tcp;
      }
    }
    congestion-notification-profile fcoe-cnp;
  }
  ae1 {
    forwarding-class-set {
      fcoe-pg {
        output-traffic-control-profile fcoe-tcp;
      }
    }
    congestion-notification-profile fcoe-cnp;
  }
}
scheduler-maps {
  fcoe-map {
    forwarding-class fcoe scheduler fcoe-sched;
  }
}
schedulers {
  fcoe-sched {
    transmit-rate 3g;
    shaping-rate percent 100;
    priority low;
  }
}
}

```

NOTE: The forwarding class and classifier configurations are not shown because the `show` command does not display default portions of the configuration.

Display the results of the CoS configuration on FCoE Transit Switch TS1 and on FCoE Transit Switch TS2 (the results on both transit switches are the same).

```

user@switch> show configuration class-of-service
traffic-control-profiles {
    fcoe-tcp {
        scheduler-map fcoe-map;
        shaping-rate percent 100;
        guaranteed-rate 3g;
    }
}
forwarding-class-sets {
    fcoe-pg {
        class fcoe;
    }
}
congestion-notification-profile {
    fcoe-cnp {
        input {
            ieee-802.1 {
                code-point 011 {
                    pfc;
                }
            }
        }
    }
}
interfaces {
    xe-0/0/30 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        congestion-notification-profile fcoe-cnp;
    }
    xe-0/0/31 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
    }
}

```

```

        congestion-notification-profile fcoe-cnp;
    }
    xe-0/0/32 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        congestion-notification-profile fcoe-cnp;
    }
    xe-0/0/33 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        congestion-notification-profile fcoe-cnp;
    }
    ae1 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        congestion-notification-profile fcoe-cnp;
    }
}
scheduler-maps {
    fcoe-map {
        forwarding-class fcoe scheduler fcoe-sched;
    }
}
schedulers {
    fcoe-sched {
        transmit-rate 3g;
        shaping-rate percent 100;
        priority low;
    }
}
}

```

Verification

IN THIS SECTION

- [Verifying That the Output Queue Schedulers Have Been Created | 559](#)
- [Verifying That the Priority Group Output Scheduler \(Traffic Control Profile\) Has Been Created | 560](#)
- [Verifying That the Forwarding Class Set \(Priority Group\) Has Been Created | 561](#)
- [Verifying That Priority-Based Flow Control Has Been Enabled | 562](#)
- [Verifying That the Interface Class of Service Configuration Has Been Created | 563](#)
- [Verifying That the Interfaces Are Correctly Configured | 566](#)
- [Verifying That FIP Snooping Is Enabled on the FCoE VLAN on FCoE Transit Switches TS1 and TS2 Access Interfaces | 569](#)
- [Verifying That the FIP Snooping Mode Is Correct on FCoE Transit Switches TS1 and TS2 | 570](#)
- [Verifying That IGMP Snooping Is Disabled on the FCoE VLAN | 571](#)

To verify that the CoS components and FIP snooping have been configured and are operating properly, perform these tasks. Because this example uses the default fcoe forwarding class and the default IEEE 802.1p trusted classifier, the verification of those configurations is not shown.

Verifying That the Output Queue Schedulers Have Been Created

Purpose

Verify that the output queue scheduler for FCoE traffic has the correct bandwidth parameters and priorities, and is mapped to the correct forwarding class (output queue). Queue scheduler verification is the same on each of the four switches.

Action

List the scheduler map using the operational mode command `show class-of-service scheduler-map fcoe-map`:

```
user@switch> show class-of-service scheduler-map fcoe-map
Scheduler map: fcoe-map, Index: 9023

Scheduler: fcoe-sched, Forwarding class: fcoe, Index: 37289
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       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 fcoe-map` command lists the properties of the scheduler map `fcoe-map`. The command output includes:

- The name of the scheduler map (`fcoe-map`)
- The name of the scheduler (`fcoe-sched`)
- The forwarding classes mapped to the scheduler (`fcoe`)
- The minimum guaranteed queue bandwidth (transmit rate 3000000000 bps)
- The scheduling priority (`low`)
- The maximum bandwidth in the priority group the queue can consume (shaping rate 100 percent)
- The drop profile loss priority for each drop profile name. This example does not include drop profiles because you do not apply drop profiles to FCoE traffic.

Verifying That the Priority Group Output Scheduler (Traffic Control Profile) Has Been Created

Purpose

Verify that the traffic control profile `fcoe-tcp` has been created with the correct bandwidth parameters and scheduler mapping. Priority group scheduler verification is the same on each of the four switches.

Action

List the FCoE traffic control profile properties using the operational mode command `show class-of-service traffic-control-profile fcoe-tcp`:

```
user@switch> show class-of-service traffic-control-profile fcoe-tcp
Traffic control profile: fcoe-tcp, Index: 18303
  Shaping rate: 100 percent
  Scheduler map: fcoe-map
  Guaranteed rate: 3000000000
```

Meaning

The `show class-of-service traffic-control-profile fcoe-tcp` 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 (`fcoe-tcp`)
- The maximum port bandwidth the priority group can consume (shaping rate 100 percent)
- The scheduler map associated with the traffic control profile (`fcoe-map`)
- The minimum guaranteed priority group port bandwidth (guaranteed rate 3000000000 in bps)

Verifying That the Forwarding Class Set (Priority Group) Has Been Created

Purpose

Verify that the FCoE priority group has been created and that the `fcoe` priority (forwarding class) belongs to the FCoE priority group. Forwarding class set verification is the same on each of the four switches.

Action

List the forwarding class sets using the operational mode command `show class-of-service forwarding-class-set fcoe-pg`:

```
user@switch> show class-of-service forwarding-class-set fcoe-pg
Forwarding class set: fcoe-pg, Type: normal-type, Forwarding class set index: 31420
  Forwarding class      Index
  fcoe                  1
```

Meaning

The `show class-of-service forwarding-class-set fcoe-pg` command lists all of the forwarding classes (priorities) that belong to the `fcoe-pg` priority group, and the internal index number of the priority group. The command output shows that the forwarding class set `fcoe-pg` includes the forwarding class `fcoe`.

Verifying That Priority-Based Flow Control Has Been Enabled

Purpose

Verify that PFC is enabled on the FCoE code point. PFC verification is the same on each of the four switches.

Action

List the FCoE congestion notification profile using the operational mode command `show class-of-service congestion-notification fcoe-cnp`:

```

user@switch> show class-of-service congestion-notification fcoe-cnp
Type: Input, Name: fcoe-cnp, Index: 6879
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 fcoe-cnp` command lists all of the IEEE 802.1p code points in the congestion notification profile that have PFC enabled. The command output shows that PFC is enabled on code point 011 (fcoe queue) 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 That the Interface Class of Service Configuration Has Been Created

Purpose

Verify that the CoS properties of the interfaces are correct. The verification output on MC-LAG Switches S1 and S2 differs from the output on FCoE Transit Switches TS1 and TS2.

Action

List the interface CoS configuration on MC-LAG Switches S1 and S2 using the operational mode command `show configuration class-of-service interfaces`:

```
user@switch> show configuration class-of-service interfaces
ae0 {
  forwarding-class-set {
    fcoe-pg {
      output-traffic-control-profile fcoe-tcp;
    }
  }
  congestion-notification-profile fcoe-cnp;
}

ae1 {
```



```

forwarding-class-set {
    fcoe-pg {
        output-traffic-control-profile fcoe-tcp;
    }
}
congestion-notification-profile fcoe-cnp;
}

```

List the interface CoS configuration on FCoE Transit Switches TS1 and TS2 using the operational mode command `show configuration class-of-service interfaces`:

```

user@switch> show configuration class-of-service interfaces
xe-0/0/30 {
    forwarding-class-set {
        fcoe-pg {
            output-traffic-control-profile fcoe-tcp;
        }
    }
    congestion-notification-profile fcoe-cnp;
}
xe-0/0/31 {
    forwarding-class-set {
        fcoe-pg {
            output-traffic-control-profile fcoe-tcp;
        }
    }
    congestion-notification-profile fcoe-cnp;
}
xe-0/0/32 {
    forwarding-class-set {
        fcoe-pg {
            output-traffic-control-profile fcoe-tcp;
        }
    }
    congestion-notification-profile fcoe-cnp;
}
xe-0/0/33 {
    forwarding-class-set {
        fcoe-pg {
            output-traffic-control-profile fcoe-tcp;
        }
    }
}

```

```

    }
    congestion-notification-profile fcoe-cnp;
}
ae1 {
    forwarding-class-set {
        fcoe-pg {
            output-traffic-control-profile fcoe-tcp;
        }
    }
    congestion-notification-profile fcoe-cnp;
}

```

Meaning

The `show configuration class-of-service interfaces` command lists the class of service configuration for all interfaces. For each interface, the command output includes:

- The name of the interface (for example, `ae0` or `xe-0/0/30`)
- The name of the forwarding class set associated with the interface (`fcoe-pg`)
- The name of the traffic control profile associated with the interface (output traffic control profile, `fcoe-tcp`)
- The name of the congestion notification profile associated with the interface (`fcoe-cnp`)

NOTE: Interfaces that are members of a LAG are not shown individually. The LAG or MC-LAG CoS configuration is applied to all interfaces that are members of the LAG or MC-LAG. For example, the interface CoS configuration output on MC-LAG Switches S1 and S2 shows the LAG CoS configuration but does not show the CoS configuration of the member interfaces separately. The interface CoS configuration output on FCoE Transit Switches TS1 and TS2 shows the LAG CoS configuration but also shows the configuration for interfaces `xe-0/0/30`, `xe-0/0/31`, `xe-0/0/32`, and `xe-0/0/33`, which are not members of a LAG.

Verifying That the Interfaces Are Correctly Configured

Purpose

Verify that the LAG membership, MTU, VLAN membership, and port mode of the interfaces are correct. The verification output on MC-LAG Switches S1 and S2 differs from the output on FCoE Transit Switches TS1 and TS2.

Action

List the interface configuration on MC-LAG Switches S1 and S2 using the operational mode command `show configuration interfaces`:

```
user@switch> show configuration interfaces
xe-0/0/10 {
    ether-options {
        802.3ad ae0;
    }
}
xe-0/0/11 {
    ether-options {
        802.3ad ae0;
    }
}
xe-0/0/20 {
    ether-options {
        802.3ad ae1;
    }
}
xe-0/0/21 {
    ether-options {
        802.3ad ae1;
    }
}
ae0 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            port-mode trunk;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}
```

```

    }
  }
}
ae1 {
  mtu 2180;
  unit 0 {
    family ethernet-switching {
      port-mode trunk;
      vlan {
        members fcoe_vlan;
      }
    }
  }
}

```

List the interface configuration on FCoE Transit Switches TS1 and TS2 using the operational mode command `show configuration interfaces`:

```

user@switch> show configuration interfaces
xe-0/0/25 {
  ether-options {
    802.3ad ae1;
  }
}
xe-0/0/26 {
  ether-options {
    802.3ad ae1;
  }
}
xe-0/0/30 {
  mtu 2180;
  unit 0 {
    family ethernet-switching {
      port-mode tagged-access;
      vlan {
        members fcoe_vlan;
      }
    }
  }
}
xe-0/0/31 {

```

```

mtu 2180;
unit 0 {
    family ethernet-switching {
        port-mode tagged-access;
        vlan {
            members fcoe_vlan;
        }
    }
}
xe-0/0/32 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            port-mode tagged-access;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}
xe-0/0/33 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            port-mode tagged-access;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}

ae1 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            port-mode trunk;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}

```

```
}
```

Meaning

The `show configuration interfaces` command lists the configuration of each interface by interface name.

For each interface that is a member of a LAG, the command lists only the name of the LAG to which the interface belongs.

For each LAG interface and for each interface that is not a member of a LAG, the command output includes:

- The MTU (2180)
- The unit number of the interface (0)
- The port mode (trunk mode for interfaces that connect two switches, tagged-access mode for interfaces that connect to FCoE hosts)
- The name of the VLAN in which the interface is a member (fcoe_vlan)

Verifying That FIP Snooping Is Enabled on the FCoE VLAN on FCoE Transit Switches TS1 and TS2 Access Interfaces

Purpose

Verify that FIP snooping is enabled on the FCoE VLAN access interfaces. FIP snooping is enabled only on the FCoE access interfaces, so it is enabled only on FCoE Transit Switches TS1 and TS2. FIP snooping is not enabled on MC-LAG Switches S1 and S2 because FIP snooping is done at the Transit Switch TS1 and TS2 FCoE access ports.

Action

List the port security configuration on FCoE Transit Switches TS1 and TS2 using the operational mode command `show configuration ethernet-switching-options secure-access-port`:

```
user@switch> show configuration ethernet-switching-options secure-access-port
interface ae1.0 {
    fcoe-trusted;
}
vlan fcoe_vlan {
    examine-fip {
```

```

        examine-vn2vn {
            beacon-period 90000;
        }
    }
}

```

Meaning

The `show configuration ethernet-switching-options secure-access-port` command lists port security information, including whether a port is trusted. The command output shows that:

- LAG port `ae1.0`, which connects the FCoE transit switch to the MC-LAG switches, is configured as an FCoE trusted interface. FIP snooping is not performed on the member interfaces of the LAG (`xe-0/0/25` and `xe-0/0/26`).
- FIP snooping is enabled (`examine-fip`) on the FCoE VLAN (`fcoe_vlan`), the type of FIP snooping is VN2VN_Port FIP snooping (`examine-vn2vn`), and the beacon period is set to 90000 milliseconds. On Transit Switches TS1 and TS2, all interface members of the FCoE VLAN perform FIP snooping unless the interface is configured as FCoE trusted. On Transit Switches TS1 and TS2, interfaces `xe-0/0/30`, `xe-0/0/31`, `xe-0/0/32`, and `xe-0/0/33` perform FIP snooping because they are not configured as FCoE trusted. The interface members of LAG `ae1` (`xe-0/0/25` and `xe-0/0/26`) do not perform FIP snooping because the LAG is configured as FCoE trusted.

Verifying That the FIP Snooping Mode Is Correct on FCoE Transit Switches TS1 and TS2

Purpose

Verify that the FIP snooping mode is correct on the FCoE VLAN. FIP snooping is enabled only on the FCoE access interfaces, so it is enabled only on FCoE Transit Switches TS1 and TS2. FIP snooping is not enabled on MC-LAG Switches S1 and S2 because FIP snooping is done at the Transit Switch TS1 and TS2 FCoE access ports.

Action

List the FIP snooping configuration on FCoE Transit Switches TS1 and TS2 using the operational mode command `show fip snooping brief`:

```

user@switch> show fip snooping brief
VLAN: fcoe_vlan,      Mode: VN2VN Snooping
FC-MAP: 0e:fd:00

```

...

NOTE: The output has been truncated to show only the relevant information.

Meaning

The `show fip snooping brief` command lists FIP snooping information, including the FIP snooping VLAN and the FIP snooping mode. The command output shows that:

- The VLAN on which FIP snooping is enabled is `fcoe_vlan`
- The FIP snooping mode is `VN2VN_Port FIP snooping (VN2VN Snooping)`

Verifying That IGMP Snooping Is Disabled on the FCoE VLAN

Purpose

Verify that IGMP snooping is disabled on the FCoE VLAN on all four switches.

Action

List the IGMP snooping protocol information on each of the four switches using the `show configuration protocols igmp-snooping` command:

```
user@switch> show configuration protocols igmp-snooping
vlan fcoe_vlan {
    disable;
}
```

Meaning

The `show configuration protocols igmp-snooping` command lists the IGMP snooping configuration for the VLANs configured on the switch. The command output shows that IGMP snooping is disabled on the FCoE VLAN (`fcoe_vlan`).

RELATED DOCUMENTATION

[Example: Configuring CoS PFC for FCoE Traffic | 524](#)

Example: Configuring CoS Using ELS for FCoE Transit Switch Traffic Across an MC-LAG

IN THIS SECTION

- [Requirements | 573](#)
- [Overview | 573](#)
- [Configuration | 580](#)
- [Verification | 595](#)

Multichassis link aggregation groups (MC-LAGs) provide redundancy and load balancing between two QFX Series switches, multihoming support for client devices such as servers, and a loop-free Layer 2 network without running Spanning Tree Protocol (STP).

NOTE: This example uses the Junos OS Enhanced Layer 2 Software (ELS) configuration style for QFX Series switches. If your switch runs software that does not support ELS, see [Example: Configuring CoS for FCoE Transit Switch Traffic Across an MC-LAG](#). For ELS details, see [Using the Enhanced Layer 2 Software CLI](#).

You can use an MC-LAG to provide a redundant aggregation layer for Fibre Channel over Ethernet (FCoE) traffic in an *inverted-U* topology. To support lossless transport of FCoE traffic across an MC-LAG, you must configure the appropriate class of service (CoS) on both of the QFX Series switches with MC-LAG port members. The CoS configuration must be the same on both of the MC-LAG switches because an MC-LAG does not carry forwarding class and IEEE 802.1p priority information.

Ports that are members of an MC-LAG act as FCoE passthrough transit switch ports.

NOTE: This example describes how to configure CoS to provide lossless transport for FCoE traffic across an MC-LAG that connects two QFX Series switches. It also describes how to

configure CoS on the FCoE transit switches that connect FCoE hosts to the QFX Series switches that form the MC-LAG.

This example does not describe how to configure the MC-LAG itself; it includes a subset of MC-LAG configuration that only shows how to configure interface membership in the MC-LAG.

This example does *not* describe how to configure the MC-LAG itself. For a detailed example of MC-LAG configuration, see [Example: Configuring Multichassis Link Aggregation on the QFX Series](#). However, this example includes a subset of MC-LAG configuration that only shows how to configure interface membership in the MC-LAG.

NOTE: Juniper Networks QFX10000 aggregation switches do not support FIP snooping, so they cannot be used as FIP snooping access switches (Transit Switches TS1 and TS2) in this example. However, QFX10000 switches can play the role of the MC-LAG switches (MC-LAG Switch S1 and MC-LAG Switch S2) in this example.

QFX3500 and QFX3600 Virtual Chassis switches do not support FCoE.

This topic describes:

Requirements

This example uses the following hardware and software components:

- Two Juniper Networks QFX5100 Switches running the ELS CLI that form an MC-LAG for FCoE traffic.
- Two Juniper Networks QFX5100 Switches running the ELS CLI that provide FCoE server access in transit switch mode and that connect to the MC-LAG switches.
- FCoE servers (or other FCoE hosts) connected to the transit switches.
- Junos OS Release 13.2 or later for the QFX Series.

Overview

IN THIS SECTION

- [Topology](#) | 574

FCoE traffic requires lossless transport. This example shows you how to:

- Configure CoS for FCoE traffic on the two QFX5100 switches that form the MC-LAG, including priority-based flow control (PFC). The example also includes configuration for both enhanced transmission selection (ETS) hierarchical scheduling of resources for the FCoE forwarding class priority and for the forwarding class set priority group, and also direct port scheduling. You can only use one of the scheduling methods on a port. Different switches support different scheduling methods.

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.

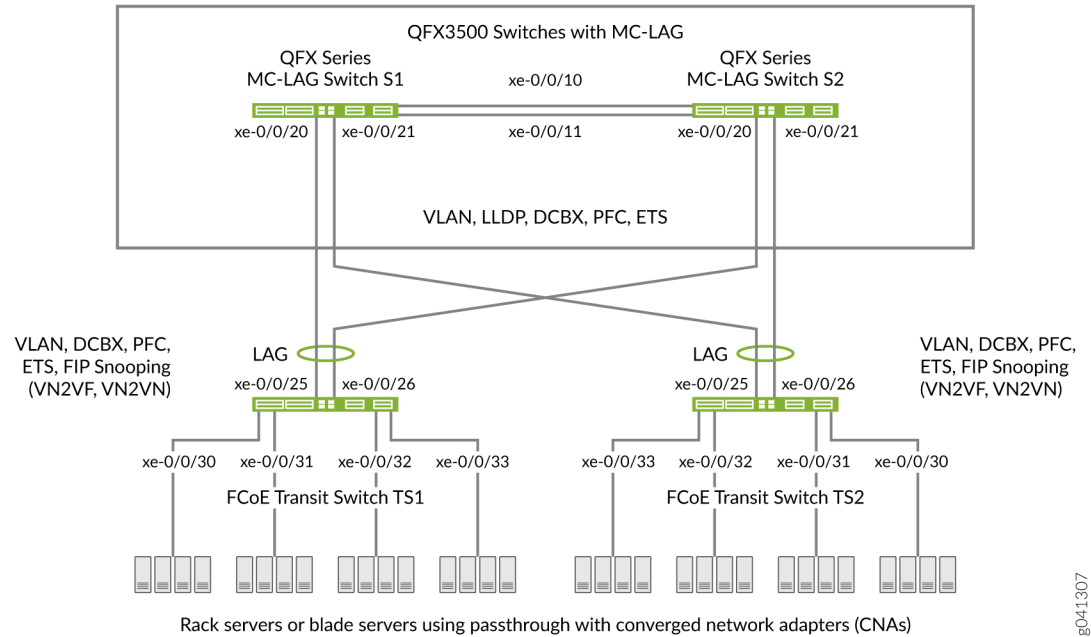
- Configure CoS for FCoE on the two FCoE transit switches that connect FCoE hosts to the MC-LAG switches and enable FIP snooping on the FCoE VLAN at the FCoE transit switch access ports.
- Configure the appropriate port mode, MTU, and FCoE trusted or untrusted state for each interface to support lossless FCoE transport.

NOTE: Do not enable IGMP snooping on the FCoE VLAN. (IGMP snooping is enabled on the default VLAN by default, but is disabled by default on all other VLANs.)

Topology

QFX5100 switches that act as transit switches support MC-LAGs for FCoE traffic in an inverted-U network topology, as shown in [Figure 25 on page 575](#).

Figure 25: Supported Topology for an MC-LAG on an FCoE Transit Switch



NOTE: Juniper Networks QFX10000 aggregation switches do not support FIP snooping, so they cannot be used as FIP snooping access switches (Transit Switches TS1 and TS2) in this example. However, QFX10000 switches can play the role of the MC-LAG switches (MC-LAG Switch S1 and MC-LAG Switch S2) in this example.

Table 91 on page 575 shows the configuration components for this example.

Table 91: Components of the CoS for FCoE Traffic Across an MC-LAG Configuration Topology

Component	Settings
Hardware	Four QFX5100 switches running the ELS CLI (two to form the MC-LAG as passthrough transit switches and two transit switches for FCoE access).
Forwarding class (all switches)	Default fcoe forwarding class.

Table 91: Components of the CoS for FCoE Traffic Across an MC-LAG Configuration Topology
(Continued)

Component	Settings
Classifier (forwarding class mapping of incoming traffic to IEEE priority)	Default IEEE 802.1p trusted classifier on all FCoE interfaces.
LAGs and MC-LAG	<p>S1—Ports xe-0/0/10 and x-0/0/11 are members of LAG ae0, which connects Switch S1 to Switch S2. Ports xe-0/0/20 and xe-0/0/21 are members of MC-LAG ae1.</p> <p>All ports are configured in trunk interface mode, as fcoe-trusted, and with an MTU of 2180.</p> <p>S2—Ports xe-0/0/10 and x-0/0/11 are members of LAG ae0, which connects Switch S2 to Switch S1. Ports xe-0/0/20 and xe-0/0/21 are members of MC-LAG ae1.</p> <p>All ports are configured in trunk interface mode, as fcoe-trusted, and with an MTU of 2180.</p> <p>NOTE: Ports xe-0/0/20 and xe-0/0/21 on Switches S1 and S2 are the members of the MC-LAG.</p> <p>TS1—Ports xe-0/0/25 and x-0/0/26 are members of LAG ae1, configured in trunk interface mode, as fcoe-trusted, and with an MTU of 2180.</p> <p>Ports xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 are configured in trunk interface mode, with an MTU of 2180.</p> <p>TS2—Ports xe-0/0/25 and x-0/0/26 are members of LAG ae1, configured in trunk interface mode, as fcoe-trusted, and with an MTU of 2180.</p> <p>Ports xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 are configured in trunk interface mode, with an MTU of 2180.</p>
FCoE queue scheduler (all switches)	<p>fcoe-sched:</p> <p>Minimum bandwidth 3g</p> <p>Maximum bandwidth 100%</p> <p>Priority low</p>

Table 91: Components of the CoS for FCoE Traffic Across an MC-LAG Configuration Topology
(Continued)

Component	Settings
Forwarding class-to-scheduler mapping (all switches)	<p>Scheduler map fcoe-map:</p> <p>Forwarding class fcoe</p> <p>Scheduler fcoe-sched</p>
PFC congestion notification profile (all switches)	<p>fcoe-cnp:</p> <p>Code point 011</p> <p>Ingress interfaces:</p> <ul style="list-style-type: none"> • S1—LAG ae0 and MC-LAG ae1 • S2—LAG ae0 and MC-LAG ae1 • TS1—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 • TS2—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33
FCoE VLAN name and tag ID	<p>Name—fcoe_vlan</p> <p>ID—100</p> <p>Include the FCoE VLAN on the interfaces that carry FCoE traffic on all four switches.</p>
ETS only—forwarding class set (FCoE priority group, all switches)	<p>fcoe-pg:</p> <p>Forwarding class fcoe</p> <p>Egress interfaces:</p> <ul style="list-style-type: none"> • S1—LAG ae0 and MC-LAG ae1 • S2—LAG ae0 and MC-LAG ae1 • TS1—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 • TS2—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33

Table 91: Components of the CoS for FCoE Traffic Across an MC-LAG Configuration Topology
(Continued)

Component	Settings
ETS only—traffic control profile (all switches)	<p>fcoe-tcp: Scheduler map fcoe-map Minimum bandwidth 3g Maximum bandwidth 100%</p> <p>The traffic control profile is applied to the same interfaces as the forwarding class set, using the same CLI statement. This applies ETS hierarchical scheduling to the interfaces.</p>
Port scheduling only—apply scheduling to interfaces	<p>On switches that support direct port scheduling, if you use port scheduling, apply scheduling by attaching the scheduler map directly to interfaces:</p> <ul style="list-style-type: none"> • S1—LAG ae0 and MC-LAG ae1 • S2—LAG ae0 and MC-LAG ae1 • TS1—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 • TS2—LAG ae1, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33
FIP snooping	<p>Enable FIP snooping on Transit Switches TS1 and TS2 on the FCoE VLAN. Configure the LAG interfaces that connect to the MC-LAG switches as FCoE trusted interfaces so that they do not perform FIP snooping.</p> <p>This example enables VN2VN_Port FIP snooping on the FCoE transit switch interfaces connected to the FCoE servers. The example is equally valid with VN2VF_Port FIP snooping enabled on the transit switch access ports. The method of FIP snooping you enable depends on your network configuration.</p> <p>NOTE: Juniper Networks QFX10000 aggregation switches do not support FIP snooping, so they cannot be used as FIP snooping access switches (Transit Switches TS1 and TS2) in this example.</p>

NOTE: This example uses the default IEEE 802.1p trusted BA classifier, which is automatically applied to trunk mode interfaces if you do not apply an explicitly configured classifier.

To configure CoS for FCoE traffic across an MC-LAG:

- Use the default FCoE forwarding class and forwarding-class-to-queue mapping (do not explicitly configure the FCoE forwarding class or output queue). The default FCoE forwarding class is `fcoe`, and the default output queue is `queue 3`.
- Use the default trusted BA classifier, which maps incoming packets to forwarding classes by the IEEE 802.1p code point (CoS priority) of the packet. The trusted classifier is the default classifier for interfaces in trunk interface mode. The default trusted classifier maps incoming packets with the IEEE 802.1p code point 3 (011) to the FCoE forwarding class. If you choose to configure the BA classifier instead of using the default classifier, you must ensure that FCoE traffic is classified into forwarding classes in exactly the same way on both MC-LAG switches. Using the default classifier ensures consistent classifier configuration on the MC-LAG ports.
- Configure a congestion notification profile that enables PFC on the FCoE code point (code point 011 in this example). The congestion notification profile configuration must be the same on both MC-LAG switches.
- Apply the congestion notification profile to the interfaces.
- Configure the interface mode, MTU, and FCoE trusted or untrusted state for each interface to support lossless FCoE transport.
- For ETS hierarchical port scheduling, configure ETS on the interfaces to provide the bandwidth required for lossless FCoE transport. Configuring ETS includes configuring bandwidth scheduling for the FCoE forwarding class, a forwarding class set (priority group) that includes the FCoE forwarding class, and a traffic control profile to assign bandwidth to the forwarding class set that includes FCoE traffic, and applying the traffic control profile and forwarding class set to interfaces..

On switches that support direct port scheduling, configure CoS properties on interfaces by applying scheduler maps directly to interfaces.

In addition, this example describes how to enable FIP snooping on the Transit Switch TS1 and TS2 ports that are connected to the FCoE servers. To provide secure access, FIP snooping must be enabled on the FCoE access ports.

This example focuses on the CoS configuration to support lossless FCoE transport across an MC-LAG. This example does not describe how to configure the properties of MC-LAGs and LAGs, although it does show you how to configure the port characteristics required to support lossless transport and how to assign interfaces to the MC-LAG and to the LAGs.

Before you configure CoS, configure:

- The MC-LAGs that connect Switches S1 and S2 to Switches TS1 and TS2. ([Example: Configuring Multichassis Link Aggregation on the QFX Series](#) describes how to configure MC-LAGs.)
- The LAGs that connect the Transit Switches TS1 and TS2 to MC-LAG Switches S1 and S2. ([Configuring Link Aggregation](#) describes how to configure LAGs.)
- The LAG that connects Switch S1 to Switch S2.

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 580](#)
- [MC-LAG Switches S1 and S2 Common Configuration \(Applies to ETS and Port Scheduling\) | 583](#)
- [MC-LAG Switches S1 and S2 ETS Hierarchical Scheduling Configuration | 585](#)
- [MC-LAG Switches S1 and S2 Port Scheduling Configuration | 586](#)
- [FCoE Transit Switches TS1 and TS2 Common Configuration \(Applies to ETS and Port Scheduling\) | 587](#)
- [FCoE Transit Switches TS1 and TS2 ETS Hierarchical Scheduling Configuration | 590](#)
- [FCoE Transit Switches TS1 and TS2 Port Scheduling Configuration | 591](#)
- [Results | 591](#)

To configure CoS for lossless FCoE transport across an MC-LAG, perform these tasks:

CLI Quick Configuration

To quickly configure CoS for lossless FCoE transport across an MC-LAG, 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 the MC-LAG and FCoE transit switches at the [edit] hierarchy level.

The quick configuration shows the commands for the two MC-LAG switches and the two FCoE transit switches separately. The configurations on both of the MC-LAG switches are same and on both of the FCoE transit switches are the same because the CoS configuration must be identical, and because this example uses the same ports on each of these sets of switches.

NOTE: The CLI configurations for the MC-LAG switches and for the FCoE transit switches are each separated into three sections:

- Configuration common to all port scheduling methods
- Configuration specific to ETS hierarchical port scheduling
- Configuration specific to direct port scheduling

Quick configuration for MC-LAG Switch S1 and Switch S2:

MC-LAG Switches Configuration Common to ETS Hierarchical Port Scheduling and to Direct Port Scheduling

```
set class-of-service schedulers fcoe-sched priority low transmit-rate 3g
set class-of-service schedulers fcoe-sched shaping-rate percent 100
set class-of-service scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
set class-of-service congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
set class-of-service interfaces ae0 congestion-notification-profile fcoe-cnp
set class-of-service interfaces ae1 congestion-notification-profile fcoe-cnp
set vlans fcoe_vlan vlan-id 100
set interfaces xe-0/0/10 ether-options 802.3ad ae0
set interfaces xe-0/0/11 ether-options 802.3ad ae0
set interfaces xe-0/0/20 ether-options 802.3ad ae1
set interfaces xe-0/0/21 ether-options 802.3ad ae1
set interfaces ae0 unit 0 family ethernet-switching interface-mode trunk vlan members fcoe_vlan
set interfaces ae1 unit 0 family ethernet-switching interface-mode trunk vlan members fcoe_vlan
set interfaces ae0 mtu 2180
set interfaces ae1 mtu 2180
set vlans fcoe_vlan forwarding-options fip-security interface ae0 fcoe-trusted
set vlans fcoe_vlan forwarding-options fip-security interface ae1 fcoe-trusted
```

MC-LAG Switches Configuration for ETS Hierarchical Port Scheduling

```
set class-of-service forwarding-class-sets fcoe-pg class fcoe
set class-of-service traffic-control-profiles fcoe-tcp scheduler-map fcoe-map guaranteed-rate 3g
set class-of-service traffic-control-profiles fcoe-tcp shaping-rate percent 100
set class-of-service interfaces ae0 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
```

```
set class-of-service interfaces ae1 forwarding-class-set fcoe-pg output-traffic-control-profile
fcoe-tcp
```

MC-LAG Switches Configuration for Direct Port Scheduling

```
set class-of-service interfaces ae0 scheduler-map fcoe-map
set class-of-service interfaces ae1 scheduler-map fcoe-map
```

Quick configuration for FCoE Transit Switch TS1 and Switch TS2:

FCoE Transit Switches Configuration Common to ETS Hierarchical Port Scheduling and to Direct Port Scheduling

```
set class-of-service schedulers fcoe-sched priority low transmit-rate 3g
set class-of-service schedulers fcoe-sched shaping-rate percent 100
set class-of-service scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
set class-of-service congestion-notification-profile fcoe-cnp input ieee-802.1 code-point 011 pfc
set class-of-service interfaces ae1 congestion-notification-profile fcoe-cnp
set class-of-service interfaces xe-0/0/30 congestion-notification-profile fcoe-cnp
set class-of-service interfaces xe-0/0/31 congestion-notification-profile fcoe-cnp
set class-of-service interfaces xe-0/0/32 congestion-notification-profile fcoe-cnp
set class-of-service interfaces xe-0/0/33 congestion-notification-profile fcoe-cnp
set vlans fcoe_vlan vlan-id 100
set interfaces xe-0/0/25 ether-options 802.3ad ae1
set interfaces xe-0/0/26 ether-options 802.3ad ae1
set interfaces ae1 unit 0 family ethernet-switching interface-mode trunk vlan members fcoe_vlan
set interfaces xe-0/0/30 unit 0 family ethernet-switching interface-mode trunk vlan members
fcoe_vlan
set interfaces xe-0/0/31 unit 0 family ethernet-switching interface-mode trunk vlan members
fcoe_vlan
set interfaces xe-0/0/32 unit 0 family ethernet-switching interface-mode trunk vlan members
fcoe_vlan
set interfaces xe-0/0/33 unit 0 family ethernet-switching interface-mode trunk vlan members
fcoe_vlan
set interfaces ae1 mtu 2180
set interfaces xe-0/0/30 mtu 2180
set interfaces xe-0/0/31 mtu 2180
set interfaces xe-0/0/32 mtu 2180
set interfaces xe-0/0/33 mtu 2180
```

```
set vlans fcoe_vlan forwarding-options fip-security interface ae1 fcoe-trusted
set vlans fcoe_vlan forwarding-options fip-security examine-vn2v2 beacon-period 90000
```

FCoE Transit Switches Configuration for ETS Hierarchical Port Scheduling

```
set class-of-service forwarding-class-sets fcoe-pg class fcoe
set class-of-service traffic-control-profiles fcoe-tcp scheduler-map fcoe-map guaranteed-rate 3g
set class-of-service traffic-control-profiles fcoe-tcp shaping-rate percent 100
set class-of-service interfaces ae1 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service interfaces xe-0/0/30 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service interfaces xe-0/0/31 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service interfaces xe-0/0/32 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
set class-of-service interfaces xe-0/0/33 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
```

FCoE Transit Switches Configuration for Direct Port Scheduling

```
set class-of-service interfaces ae1 scheduler-map fcoe-map
set class-of-service interfaces xe-0/0/30 scheduler-map fcoe-map
set class-of-service interfaces xe-0/0/31 scheduler-map fcoe-map
set class-of-service interfaces xe-0/0/32 scheduler-map fcoe-map
set class-of-service interfaces xe-0/0/33 scheduler-map fcoe-map
```

MC-LAG Switches S1 and S2 Common Configuration (Applies to ETS and Port Scheduling)

Step-by-Step Procedure

To configure queue scheduling, PFC, the FCoE VLAN, and LAG and MC-LAG interface membership and characteristics to support lossless FCoE transport across an MC-LAG (this example uses the default `fcoe` forwarding class and the default classifier to map incoming FCoE traffic to the FCoE IEEE 802.1p code point 011), for both ETS hierarchical port scheduling and port scheduling (common configuration):

1. 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
```

2. Map the FCoE forwarding class to the FCoE scheduler (fcoe-sched):

```
[edit class-of-service]
user@switch# set scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
```

3. Enable PFC on the FCoE priority by creating a congestion notification profile (fcoe-cnp) that applies 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
```

4. Apply the PFC configuration to the LAG and MC-LAG interfaces:

```
[edit class-of-service]
user@switch# set interfaces ae0 congestion-notification-profile fcoe-cnp
user@switch# set interfaces ae1 congestion-notification-profile fcoe-cnp
```

5. Configure the VLAN for FCoE traffic (fcoe_vlan):

```
[edit vlans]
user@switch# set fcoe_vlan vlan-id 100
```

6. Add the member interfaces to the LAG between the two MC-LAG switches:

```
[edit interfaces]
user@switch# set xe-0/0/10 ether-options 802.3ad ae0
user@switch# set xe-0/0/11 ether-options 802.3ad ae0
```

7. Add the member interfaces to the MC-LAG:

```
[edit interfaces]
user@switch# set xe-0/0/20 ether-options 802.3ad ae1
user@switch# set xe-0/0/21 ether-options 802.3ad ae1
```

8. Configure the interface mode as trunk and membership in the FCoE VLAN (fcoe_vlan) for the LAG (ae0) and for the MC-LAG (ae1):

```
[edit interfaces]
user@switch# set interfaces ae0 unit 0 family ethernet-switching interface-mode trunk vlan
members fcoe_vlan
user@switch# set interfaces ae1 unit 0 family ethernet-switching interface-mode trunk vlan
members fcoe_vlan
```

9. Set the MTU to 2180 for the LAG and MC-LAG interfaces. 2180 bytes is the minimum size required to handle FCoE packets because of the payload and header sizes; you can configure the MTU to a higher number of bytes if desired, but not less than 2180 bytes:

```
[edit interfaces]
user@switch# set ae0 mtu 2180
user@switch# set ae1 mtu 2180
```

10. Set the LAG and MC-LAG interfaces as FCoE trusted ports. Ports that connect to other switches should be trusted and should not perform FIP snooping:

```
[edit]
user@switch# set vlans fcoe_vlan forwarding-options fip-security interface ae0 fcoe-trusted
user@switch# set vlans fcoe_vlan forwarding-options fip-security interface ae1 fcoe-trusted
```

MC-LAG Switches S1 and S2 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 (fcoe-pg) for the FCoE traffic:

```
[edit class-of-service]
user@switch# set forwarding-class-sets fcoe-pg class fcoe
```

2. Define the traffic control profile (fcoe-tcp) to use on 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 LAG and MC-LAG interfaces:

```
[edit class-of-service]
user@switch# set interfaces ae0 forwarding-class-set fcoe-pg output-traffic-control-profile
fcoe-tcp
user@switch# set interfaces ae1 forwarding-class-set fcoe-pg output-traffic-control-profile
fcoe-tcp
```

MC-LAG Switches S1 and S2 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:

```
set class-of-service interfaces ae0 scheduler-map fcoe-map
set class-of-service interfaces ae1 scheduler-map fcoe-map
```

FCoE Transit Switches TS1 and TS2 Common Configuration (Applies to ETS and Port Scheduling)

Step-by-Step Procedure

The CoS configuration on FCoE Transit Switches TS1 and TS2 is similar to the CoS configuration on MC-LAG Switches S1 and S2. However, the port configurations differ, and you must enable FIP snooping on the Switch TS1 and Switch TS2 FCoE access ports.

To configure queue scheduling, PFC, the FCoE VLAN, and LAG interface membership and characteristics to support lossless FCoE transport across the MC-LAG (this example uses the default `fcoe` forwarding class and the default classifier to map incoming FCoE traffic to the FCoE IEEE 802.1p code point 011, so you do not configure them), or both ETS hierarchical scheduling and port scheduling (common configuration):

1. 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
```

2. Map the FCoE forwarding class to the FCoE scheduler (`fcoe-sched`):

```
[edit class-of-service]
user@switch# set scheduler-maps fcoe-map forwarding-class fcoe scheduler fcoe-sched
```

3. Enable PFC on the FCoE priority by creating a congestion notification profile (`fcoe-cnp`) that applies 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
```

4. Apply the PFC configuration to the LAG interface and to the FCoE access interfaces:

```
[edit class-of-service]
user@switch# set interfaces ae1 congestion-notification-profile fcoe-cnp
user@switch# set class-of-service interfaces xe-0/0/30 congestion-notification-profile fcoe-cnp
```



```

user@switch# set class-of-service interfaces xe-0/0/31 congestion-notification-profile fcoe-
cnp
user@switch# set class-of-service interfaces xe-0/0/32 congestion-notification-profile fcoe-
cnp
user@switch# set class-of-service interfaces xe-0/0/33 congestion-notification-profile fcoe-
cnp

```

5. Configure the VLAN for FCoE traffic (fcoe_vlan):

```

[edit vlans]
user@switch# set fcoe_vlan vlan-id 100

```

6. Add the member interfaces to the LAG:

```

[edit interfaces]
user@switch# set xe-0/0/25 ether-options 802.3ad ae1
user@switch# set xe-0/0/26 ether-options 802.3ad ae1

```

7. On the LAG (ae1), configure the interface mode as trunk and membership in the FCoE VLAN (fcoe_vlan):

```

[edit interfaces]
user@switch# set interfaces ae1 unit 0 family ethernet-switching interface-mode trunk vlan
members fcoe_vlan

```

8. On the FCoE access interfaces (xe-0/0/30, xe-0/0/31, xe-0/0/32, xe-0/0/33), configure the interface mode as trunk and membership in the FCoE VLAN (fcoe_vlan):

```

[edit interfaces]
user@switch# set interfaces xe-0/0/30 unit 0 family ethernet-switching interface-mode trunk
vlan members fcoe_vlan
user@switch# set interfaces xe-0/0/31 unit 0 family ethernet-switching interface-mode trunk
vlan members fcoe_vlan
user@switch# set interfaces xe-0/0/32 unit 0 family ethernet-switching interface-mode trunk
vlan members fcoe_vlan

```

```
user@switch# set interfaces xe-0/0/33 unit 0 family ethernet-switching interface-mode trunk
vlan members fcoe_vlan
```

9. Set the MTU to 2180 for the LAG and FCoE access interfaces. 2180 bytes is the minimum size required to handle FCoE packets because of the payload and header sizes; you can configure the MTU to a higher number of bytes if desired, but not less than 2180 bytes:

```
[edit interfaces]
user@switch# set ae1 mtu 2180
user@switch# set xe-0/0/30 mtu 2180
user@switch# set xe-0/0/31 mtu 2180
user@switch# set xe-0/0/32 mtu 2180
user@switch# set xe-0/0/33 mtu 2180
```

10. Set the LAG interface as an FCoE trusted port. Ports that connect to other switches should be trusted and should not perform FIP snooping:

```
[edit]
user@switch# set vlans fcoe_vlan forwarding-options fip-security interface ae1 fcoe-trusted
```

NOTE: Access ports xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 are not configured as FCoE trusted ports. The access ports remain in the default state as untrusted ports because they connect directly to FCoE devices and must perform FIP snooping to ensure network security.

11. Enable FIP snooping on the FCoE VLAN to prevent unauthorized FCoE network access (this example uses VN2VN_Port FIP snooping; the example is equally valid if you use VN2VF_Port FIP snooping):

```
[edit]
user@switch# set vlans fcoe_vlan forwarding-options fip-security examine-vn2vn beacon-
period 90000
```

NOTE: QFX10000 switches do not support FIP snooping and cannot be used as FCoE access transit switches. (QFX10000 switches can be used as FCoE aggregation switches.)

FCoE Transit Switches TS1 and TS2 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 (fcoe-pg) for the FCoE traffic:

```
[edit class-of-service]
user@switch# set forwarding-class-sets fcoe-pg class fcoe
```

2. Define the traffic control profile (fcoe-tcp) to use on 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 LAG interface and to the FCoE access interfaces:

```
[edit class-of-service]
user@switch# set interfaces ae1 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
user@switch# set class-of-service interfaces xe-0/0/30 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
user@switch# set class-of-service interfaces xe-0/0/31 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
user@switch# set class-of-service interfaces xe-0/0/32 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
user@switch# set class-of-service interfaces xe-0/0/33 forwarding-class-set fcoe-pg output-traffic-control-profile fcoe-tcp
```

FCoE Transit Switches TS1 and TS2 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:

```
user@switch# set class-of-service interfaces ae1 scheduler-map fcoe-map
user@switch# set class-of-service interfaces xe-0/0/30 scheduler-map fcoe-map
user@switch# set class-of-service interfaces xe-0/0/31 scheduler-map fcoe-map
user@switch# set class-of-service interfaces xe-0/0/32 scheduler-map fcoe-map
user@switch# set class-of-service interfaces xe-0/0/33 scheduler-map fcoe-map
```

Results

Display the results of the CoS configuration on MC-LAG Switch S1 and on MC-LAG Switch S2 (the results on both switches are the same). The results are from the ETS hierarchical scheduling configuration, which shows the more complex configuration. Direct port scheduling results would not show the traffic control profile or forwarding class set portions of the configuration, but 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). Other than that, they are the same.

```
user@switch> show configuration class-of-service
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 {
```


For MC-LAG verification commands, see [Example: Configuring Multichassis Link Aggregation on the QFX Series](#).

Display the results of the CoS configuration on FCoE Transit Switch TS1 and on FCoE Transit Switch TS2 (the results on both transit switches are the same). The results are from the ETS hierarchical port scheduling configuration, which shows the more complex configuration. Direct port scheduling results would not show the traffic control profile or forwarding class set portions of the configuration, but 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). Other than that, they are the same.

```
user@switch> show configuration class-of-service
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-0/0/30 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
    }
}
```

```

        congestion-notification-profile fcoe-cnp;
    }
    xe-0/0/31 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        congestion-notification-profile fcoe-cnp;
    }
    xe-0/0/32 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        congestion-notification-profile fcoe-cnp;
    }
    xe-0/0/33 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        congestion-notification-profile fcoe-cnp;
    }
    ae1 {
        forwarding-class-set {
            fcoe-pg {
                output-traffic-control-profile fcoe-tcp;
            }
        }
        congestion-notification-profile fcoe-cnp;
    }
}
scheduler-maps {
    fcoe-map {
        forwarding-class fcoe scheduler fcoe-sched;
    }
}
schedulers {
    fcoe-sched {
        transmit-rate 3000000000;
    }
}

```

```

        shaping-rate percent 100;
        priority low;
    }
}

```

NOTE: The forwarding class and classifier configurations are not shown because the `show` command does not display default portions of the configuration.

Verification

IN THIS SECTION

- [Verifying That the Output Queue Schedulers Have Been Created | 595](#)
- [Verifying That the Priority Group Output Scheduler \(Traffic Control Profile\) Has Been Created \(ETS Configuration Only\) | 597](#)
- [Verifying That the Forwarding Class Set \(Priority Group\) Has Been Created \(ETS Configuration Only\) | 597](#)
- [Verifying That Priority-Based Flow Control Has Been Enabled | 598](#)
- [Verifying That the Interface Class of Service Configuration Has Been Created | 599](#)
- [Verifying That the Interfaces Are Correctly Configured | 602](#)
- [Verifying That FIP Snooping Is Enabled on the FCoE VLAN on FCoE Transit Switches TS1 and TS2 Access Interfaces | 605](#)
- [Verifying That the FIP Snooping Mode Is Correct on FCoE Transit Switches TS1 and TS2 | 606](#)

To verify that the CoS components and FIP snooping have been configured and are operating properly, perform these tasks. Because this example uses the default `fcoe` forwarding class and the default IEEE 802.1p trusted classifier, the verification of those configurations is not shown:

Verifying That the Output Queue Schedulers Have Been Created

Purpose

Verify that the output queue scheduler for FCoE traffic has the correct bandwidth parameters and priorities, and is mapped to the correct forwarding class (output queue). Queue scheduler verification is the same on each of the four switches.

Action

List the scheduler map using the operational mode command `show class-of-service scheduler-map fcoe-map`:

```
user@switch> show class-of-service scheduler-map fcoe-map
Scheduler map: fcoe-map, Index: 9023

Scheduler: fcoe-sched, Forwarding class: fcoe, Index: 37289
  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       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 fcoe-map` command lists the properties of the scheduler map `fcoe-map`. The command output includes:

- The name of the scheduler map (`fcoe-map`)
- The name of the scheduler (`fcoe-sched`)
- The forwarding classes mapped to the scheduler (`fcoe`)
- The minimum guaranteed queue bandwidth (transmit rate 3000000000 bps)
- The scheduling priority (`low`)
- The maximum bandwidth in the priority group the queue can consume (shaping rate 100 percent)
- The drop profile loss priority for each drop profile name. This example does not include drop profiles because you do not apply drop profiles to FCoE traffic.

Verifying That the Priority Group Output Scheduler (Traffic Control Profile) Has Been Created (ETS Configuration Only)

Purpose

Verify that the traffic control profile `fcoe-tcp` has been created with the correct bandwidth parameters and scheduler mapping. Priority group scheduler verification is the same on each of the four switches.

Action

List the FCoE traffic control profile properties using the operational mode command `show class-of-service traffic-control-profile fcoe-tcp`:

```
user@switch> show class-of-service traffic-control-profile fcoe-tcp
Traffic control profile: fcoe-tcp, Index: 18303
  Shaping rate: 100 percent
  Scheduler map: fcoe-map
  Guaranteed rate: 3000000000
```

Meaning

The `show class-of-service traffic-control-profile fcoe-tcp` 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 (`fcoe-tcp`)
- The maximum port bandwidth the priority group can consume (shaping rate 100 percent)
- The scheduler map associated with the traffic control profile (`fcoe-map`)
- The minimum guaranteed priority group port bandwidth (guaranteed rate 3000000000 in bps)

Verifying That the Forwarding Class Set (Priority Group) Has Been Created (ETS Configuration Only)

Purpose

Verify that the FCoE priority group has been created and that the `fcoe` priority (forwarding class) belongs to the FCoE priority group. Forwarding class set verification is the same on each of the four switches.

Action

List the forwarding class sets using the operational mode command `show class-of-service forwarding-class-set fcoe-pg`:

```
user@switch> show class-of-service forwarding-class-set fcoe-pg
Forwarding class set: fcoe-pg, Type: normal-type, Forwarding class set index: 31420
  Forwarding class      Index
  fcoe                  1
```

Meaning

The `show class-of-service forwarding-class-set fcoe-pg` command lists all of the forwarding classes (priorities) that belong to the `fcoe-pg` priority group, and the internal index number of the priority group. The command output shows that the forwarding class set `fcoe-pg` includes the forwarding class `fcoe`.

Verifying That Priority-Based Flow Control Has Been Enabled

Purpose

Verify that PFC is enabled on the FCoE code point. PFC verification is the same on each of the four switches.

Action

List the FCoE congestion notification profile using the operational mode command `show class-of-service congestion-notification fcoe-cnp`:

```
user@switch> show class-of-service congestion-notification fcoe-cnp
Type: Input, Name: fcoe-cnp, Index: 6879
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 fcoe-cnp` command lists all of the IEEE 802.1p code points in the congestion notification profile that have PFC enabled. The command output shows that PFC is enabled on code point 011 (fcoe queue) 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 That the Interface Class of Service Configuration Has Been Created

Purpose

Verify that the CoS properties of the interfaces are correct. The verification output on MC-LAG Switches S1 and S2 differs from the output on FCoE Transit Switches TS1 and TS2.

NOTE: The output is from the ETS hierarchical port scheduling configuration to show the more complex configuration. Direct port scheduling results do not show the traffic control profile or forwarding class sets because those elements are configured only for ETS. Instead, the name of the scheduler map is displayed under each interface.

Action

List the interface CoS configuration on MC-LAG Switches S1 and S2 using the operational mode command `show configuration class-of-service interfaces`:

```
user@switch> show configuration class-of-service interfaces
ae0 {
  forwarding-class-set {
    fcoe-pg {
      output-traffic-control-profile fcoe-tcp;
    }
  }
  congestion-notification-profile fcoe-cnp;
}

ae1 {
  forwarding-class-set {
    fcoe-pg {
      output-traffic-control-profile fcoe-tcp;
    }
  }
  congestion-notification-profile fcoe-cnp;
}
```

List the interface CoS configuration on FCoE Transit Switches TS1 and TS2 using the operational mode command `show configuration class-of-service interfaces`:

```
user@switch> show configuration class-of-service interfaces
xe-0/0/30 {
  forwarding-class-set {
    fcoe-pg {
      output-traffic-control-profile fcoe-tcp;
    }
  }
  congestion-notification-profile fcoe-cnp;
}
xe-0/0/31 {
  forwarding-class-set {
    fcoe-pg {
      output-traffic-control-profile fcoe-tcp;
    }
  }
}
```

```

    }
  }
  congestion-notification-profile fcoe-cnp;
}
xe-0/0/32 {
  forwarding-class-set {
    fcoe-pg {
      output-traffic-control-profile fcoe-tcp;
    }
  }
  congestion-notification-profile fcoe-cnp;
}
xe-0/0/33 {
  forwarding-class-set {
    fcoe-pg {
      output-traffic-control-profile fcoe-tcp;
    }
  }
  congestion-notification-profile fcoe-cnp;
}
ae1 {
  forwarding-class-set {
    fcoe-pg {
      output-traffic-control-profile fcoe-tcp;
    }
  }
  congestion-notification-profile fcoe-cnp;
}

```

Meaning

The `show configuration class-of-service interfaces` command lists the class of service configuration for all interfaces. For each interface, the command output includes:

- The name of the interface (for example, `ae0` or `xe-0/0/30`)
- The name of the forwarding class set associated with the interface (`fcoe-pg`)
- The name of the traffic control profile associated with the interface (output traffic control profile, `fcoe-tcp`)
- The name of the congestion notification profile associated with the interface (`fcoe-cnp`)

NOTE: Interfaces that are members of a LAG are not shown individually. The LAG or MC-LAG CoS configuration is applied to all interfaces that are members of the LAG or MC-LAG. For example, the interface CoS configuration output on MC-LAG Switches S1 and S2 shows the LAG CoS configuration but does not show the CoS configuration of the member interfaces separately. The interface CoS configuration output on FCoE Transit Switches TS1 and TS2 shows the LAG CoS configuration but also shows the configuration for interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33, which are not members of a LAG.

Verifying That the Interfaces Are Correctly Configured

Purpose

Verify that the LAG membership, MTU, VLAN membership, and port mode of the interfaces are correct. The verification output on MC-LAG Switches S1 and S2 differs from the output on FCoE Transit Switches T1 and T2.

Action

List the interface configuration on MC-LAG Switches S1 and S2 using the operational mode command `show configuration interfaces`:

```
user@switch> show configuration interfaces
xe-0/0/10 {
    ether-options {
        802.3ad ae0;
    }
}
xe-0/0/11 {
    ether-options {
        802.3ad ae0;
    }
}
xe-0/0/20 {
    ether-options {
        802.3ad ae1;
    }
}
xe-0/0/21 {
    ether-options {
```

```

        802.3ad ae1;
    }
}
ae0 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            interface-mode trunk;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}
ae1 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            interface-mode trunk;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}

```

List the interface configuration on FCoE Transit Switches TS1 and TS2 using the operational mode command `show configuration interfaces`:

```

user@switch> show configuration interfaces
xe-0/0/25 {
    ether-options {
        802.3ad ae1;
    }
}
xe-0/0/26 {
    ether-options {
        802.3ad ae1;
    }
}
xe-0/0/30 {

```



```

mtu 2180;
unit 0 {
    family ethernet-switching {
        interface-mode trunk;
        vlan {
            members fcoe_vlan;
        }
    }
}
xe-0/0/31 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            interface-mode trunk;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}
xe-0/0/32 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            interface-mode trunk;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}
xe-0/0/33 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            interface-mode trunk;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}

```

```

ae1 {
    mtu 2180;
    unit 0 {
        family ethernet-switching {
            interface-mode trunk;
            vlan {
                members fcoe_vlan;
            }
        }
    }
}

```

Meaning

The `show configuration interfaces` command lists the configuration of each interface by interface name.

For each interface that is a member of a LAG, the command lists only the name of the LAG to which the interface belongs.

For each LAG interface and for each interface that is not a member of a LAG, the command output includes:

- The MTU (2180)
- The unit number of the interface (0)
- The interface mode (trunk mode both for interfaces that connect two switches and for interfaces that connect to FCoE hosts)
- The name of the VLAN in which the interface is a member (fcoe_vlan)

Verifying That FIP Snooping Is Enabled on the FCoE VLAN on FCoE Transit Switches TS1 and TS2 Access Interfaces

Purpose

Verify that FIP snooping is enabled on the FCoE VLAN access interfaces. FIP snooping is enabled only on the FCoE access interfaces, so it is enabled only on FCoE Transit Switches TS1 and TS2. FIP snooping is not enabled on MC-LAG Switches S1 and S2 because FIP snooping is done at the Transit Switch TS1 and TS2 FCoE access ports.

Action

List the port security configuration on FCoE Transit Switches TS1 and TS2 using the operational mode command `show configuration vlans fcoe_vlan forwarding-options fip-security`:

```
user@switch> show configuration vlans fcoe_vlan forwarding-options fip-security
interface ae1.0 {
    fcoe-trusted;
}
examine-vn2vn {
    beacon-period 90000;
}
```

Meaning

The `show configuration vlans fcoe_vlan forwarding-options fip-security` command lists VLAN FIP security information, including whether a port member of the VLAN is trusted. The command output shows that:

- LAG port ae1.0, which connects the FCoE transit switch to the MC-LAG switches, is configured as an FCoE trusted interface. FIP snooping is not performed on the member interfaces of the LAG (xe-0/0/25 and xe-0/0/26).
- VN2VN_Port FIP snooping is enabled (examine-vn2vn) on the FCoE VLAN and the beacon period is set to 90000 milliseconds. On Transit Switches TS1 and TS2, all interface members of the FCoE VLAN perform FIP snooping unless the interface is configured as FCoE trusted. On Transit Switches TS1 and TS2, interfaces xe-0/0/30, xe-0/0/31, xe-0/0/32, and xe-0/0/33 perform FIP snooping because they are not configured as FCoE trusted. The interface members of LAG ae1 (xe-0/0/25 and xe-0/0/26) do not perform FIP snooping because the LAG is configured as FCoE trusted.

Verifying That the FIP Snooping Mode Is Correct on FCoE Transit Switches TS1 and TS2

Purpose

Verify that the FIP snooping mode is correct on the FCoE VLAN. FIP snooping is enabled only on the FCoE access interfaces, so it is enabled only on FCoE Transit Switches TS1 and TS2. FIP snooping is not enabled on MC-LAG Switches S1 and S2 because FIP snooping is done at the Transit Switch TS1 and TS2 FCoE access ports.

Action

List the FIP snooping configuration on FCoE Transit Switches TS1 and TS2 using the operational mode command `show fip snooping brief`:

```
user@switch> show fip snooping brief
VLAN: fcoe_vlan,    Mode: VN2VN Snooping
FC-MAP: 0e:fc:00
...
```

NOTE: The output has been truncated to show only the relevant information.

Meaning

The `show fip snooping brief` command lists FIP snooping information, including the FIP snooping VLAN and the FIP snooping mode. The command output shows that:

- The VLAN on which FIP snooping is enabled is `fcoe_vlan`
- The FIP snooping mode is `VN2VN_Port FIP snooping (VN2VN Snooping)`

RELATED DOCUMENTATION

[Example: Configuring Multichassis Link Aggregation on the QFX Series](#)

[Configuring Link Aggregation](#)

[Example: Configuring CoS PFC for FCoE Traffic | 524](#)

[Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)](#)

[Example: Configuring Queue Schedulers for Port Scheduling](#)

Understanding MC-LAGs on an FCoE Transit Switch

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 | 608](#)
- [Overview | 608](#)
- [Configuration | 611](#)
- [Verification | 614](#)

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 | 609](#)

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 92 on page 610](#) shows the configuration components for this example.

Table 92: Components of the Configuration Topology for FCoE Traffic That Does Not Use Priority 3

Component	Settings
Hardware	One switch
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 92: Components of the Configuration Topology for FCoE Traffic That Does Not Use Priority 3
(Continued)

Component	Settings
DCBX application mapping	<p>Application name—fcoe_p5_app</p> <p>Application EtherType—0x8906</p> <p>Application map name—fcoe_p5_app_map</p> <p>Application map code points—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>

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 | 611](#)
- [Configuring A Lossless FCoE Forwarding Class On IEEE 802.1p Priority 5 | 612](#)

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 | 614](#)
- [Verifying the Behavior Aggregate Classifier Configuration | 615](#)
- [Verifying the PFC Flow Control Configuration \(CNP\) | 616](#)
- [Verifying the Interface Configuration | 617](#)
- [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 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 | 633](#)

[Example: Configuring Two or More Lossless FCoE Priorities on the Same FCoE Transit Switch Interface | 620](#)

[Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications \(FCoE and iSCSI\) | 652](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

[Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 194](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

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

IN THIS SECTION

- [Requirements | 620](#)
- [Overview | 621](#)
- [Configuration | 624](#)
- [Verification | 627](#)

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 | 622](#)

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 93 on page 622](#) shows the configuration components for this example.

Table 93: 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 93: 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 | 624](#)
- [Procedure | 625](#)

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 | 627](#)
- [Verifying the Behavior Aggregate Classifier Configuration | 628](#)
- [Verifying the PFC Flow Control Configuration \(CNP\) | 629](#)
- [Verifying the Interface Configuration | 630](#)
- [Verifying the DCBX Application Configuration | 631](#)
- [Verifying the DCBX Application Map Configuration | 631](#)
- [Verifying the DCBX Application Protocol Exchange Interface Configuration | 632](#)

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 | 633](#)

[Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic \(FCoE Transit Switch\) | 608](#)

[Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications \(FCoE and iSCSI\) | 652](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

[Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 194](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

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

IN THIS SECTION

- [Requirements | 633](#)
- [Overview | 634](#)
- [Configuration | 639](#)
- [Verification | 644](#)

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 | 635](#)

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 26 on page 635 shows the topology for this example, and Table 94 on page 636 shows the configuration components for this example.

Figure 26: Topology of the Two Lossless FCoE Priorities Example

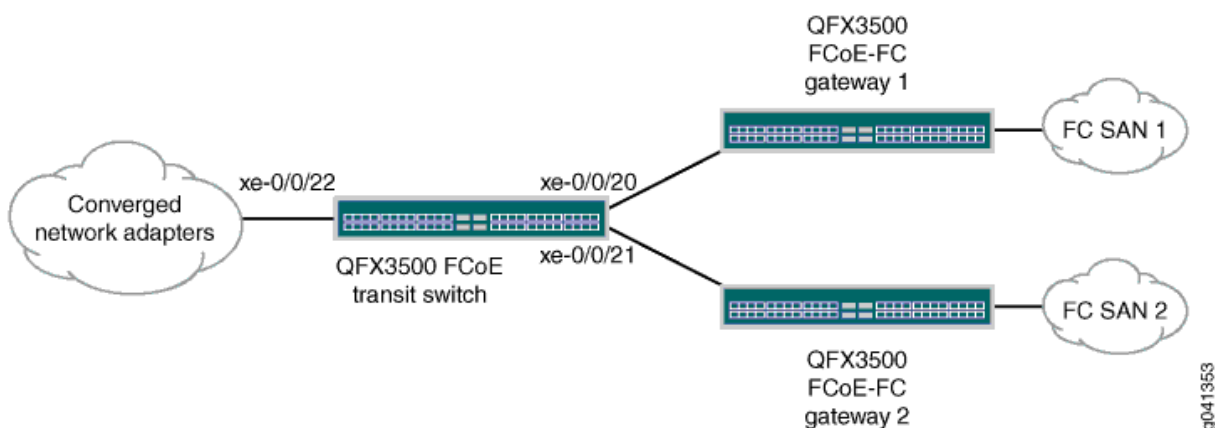


Table 94: Components of the Two Lossless FCoE Priorities Configuration Topology

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 94: 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 94: 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 | 639](#)
- [Procedure | 640](#)

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
```


- Verifying the Forwarding Class Configuration | 644
- Verifying the Behavior Aggregate Classifier Configuration | 645
- Verifying the PFC Flow Control Configuration (CNP) | 646
- Verifying the Interface Configuration | 648
- Verifying the DCBX Application Configuration | 650
- Verifying the DCBX Application Map Configuration | 650
- Verifying the DCBX Application Protocol Exchange Interface Configuration | 651

Verifying the Forwarding Class Configuration

Verify that the lossless forwarding class `fcoe1` has been created.

Show the forwarding class configuration by using the operational command `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 | 620](#)

[Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic \(FCoE Transit Switch\) | 608](#)

[Example: Configuring Lossless IEEE 802.1p Priorities on Ethernet Interfaces for Multiple Applications \(FCoE and iSCSI\) | 652](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

[Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 194](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

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

IN THIS SECTION

- [Requirements | 653](#)
- [Overview | 653](#)
- [Configuration | 660](#)
- [Verification | 668](#)

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 | 654](#)

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 27 on page 655](#) shows the topology for this example, and [Table 95 on page 655](#) shows the configuration components for this example.

Figure 27: Topology of the Lossless FCoE and iSCSI Priorities Example

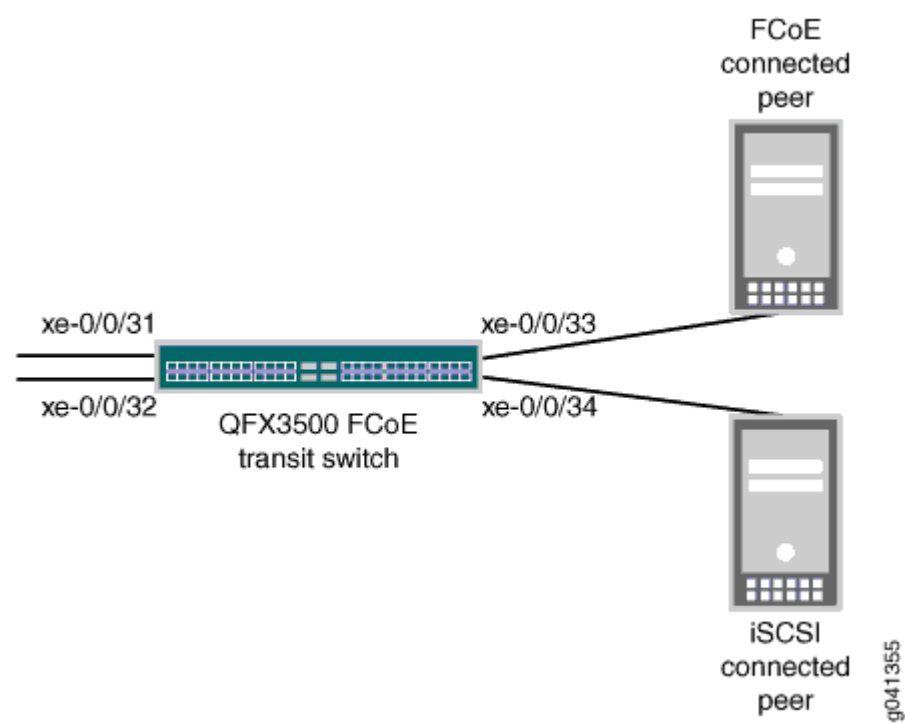


Table 95: Components of the Lossless FCoE and iSCSI Priorities Configuration Topology

Component	Settings
Hardware	One switch

Table 95: 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 95: 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 95: 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 95: 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 | 660](#)
- [Procedure | 663](#)

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 | 668](#)
- [Verifying the Behavior Aggregate Classifier Configuration | 669](#)
- [Verifying the PFC Flow Control Configuration \(CNP\) | 671](#)
- [Verifying the Interface Configuration | 674](#)
- [Verifying the DCBX Application Configuration | 676](#)
- [Verifying the DCBX Application Map Configuration | 676](#)
- [Verifying the DCBX Application Protocol Exchange Interface Configuration | 677](#)

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 | 620](#)

[Example: Configuring Lossless FCoE Traffic When the Converged Ethernet Network Does Not Use IEEE 802.1p Priority 3 for FCoE Traffic \(FCoE Transit Switch\) | 608](#)

[Example: Configuring Two or More Lossless FCoE IEEE 802.1p Priorities on Different FCoE Transit Switch Interfaces | 633](#)

[Example: Configuring DCBX Application Protocol TLV Exchange | 509](#)

[Configuring CoS PFC \(Congestion Notification Profiles\) | 216](#)

[Understanding CoS IEEE 802.1p Priorities for Lossless Traffic Flows | 194](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

Troubleshooting Dropped FCoE Traffic

IN THIS SECTION

- [Problem | 679](#)
- [Cause | 679](#)
- [Solution | 680](#)

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" on page 524](#) 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 | 524](#)

[Understanding CoS Flow Control \(Ethernet PAUSE and PFC\) | 220](#)

5

PART

CoS Buffers and the Shared Buffer Pool

[CoS Buffers Overview](#) | 684

[Shared Buffer Pool Examples](#) | 717

CoS Buffers Overview

IN THIS CHAPTER

- Understanding CoS Buffer Configuration | 684
- Configuring Global Ingress and Egress Shared Buffers | 708
- Configuring Ingress and Egress Dedicated Buffers | 710

Understanding CoS Buffer Configuration

IN THIS SECTION

- Buffer Pools | 686
- Default Buffer Pool Values | 696
- Shared Buffer Configuration Recommendations for Different Network Traffic Scenarios | 700
- Optimizing Buffer Configuration | 705
- General Buffer Configuration Rules and Considerations | 706

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 96 on page 685](#).

Table 96: 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 (`fcfs` 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 multidestination (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 97 on page 691](#):

Table 97: Default Dedicated Buffer Allocation to Egress Queues (Based on Default Scheduler)

Forwarding Class	Queue	Minimum Guaranteed Bandwidth (<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 97 on page 691](#) 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 \(ETS\)" on page 438](#) and shown in ["Example: Configuring CoS Hierarchical Port Scheduling \(ETS\)" on page 445](#). 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 98 on page 692](#) 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 99 on page 693](#) 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 98: 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%

Table 98: Egress Queue Dedicated Buffer Allocation (Example 1) (Continued)

Forwarding Class Set (Priority Group)	Forwarding Class	Queue	Scheduler Buffer Size Configuration	Buffer Allocation per Queue (Percentage)
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 99 on page 693](#). The queues with explicitly configured percentages receive the configured percentage of dedicated buffers.

Table 99: 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%

Table 99: Egress Queue Dedicated Buffer Allocation with Another Remainder Queue (Example 2)
(Continued)

Priority Group (fc-set)	Forwarding Class	Queue	Scheduler Buffer Size Configuration	Buffer Allocation per Queue (Percentage)
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 100 on page 697](#) shows the default ingress shared buffer allocation values in KB units for QFX5210 switches.

Table 100: 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 101 on page 697](#) shows the default ingress shared buffer allocation values in KB units for QFX5200-48Y switches.

Table 101: 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 102 on page 697](#) shows the default ingress shared buffer allocation values in KB units for QFX5110 and QFX5200-32C switches.

Table 102: 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 103 on page 697](#) shows the default ingress shared buffer allocation values in KB units for QFX5100, EX4600, and OCX Series switches.

Table 103: 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 104 on page 698](#) shows the default ingress shared buffer allocation values in KB units for QFX3500 and QFX3600 switches.

Table 104: 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 105 on page 698 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 105: 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 106 on page 698 shows the default egress shared buffer allocation values in KB units for QFX5210 switches.

Table 106: 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 107 on page 698 shows the default egress shared buffer allocation values in KB units for QFX5200-48Y switches.

Table 107: 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 108 on page 699 shows the default egress shared buffer allocation values in KB units for QFX5110 and QFX5200-32C switches.

Table 108: 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 109 on page 699](#) shows the default egress shared buffer allocation values in KB units for QFX5100, EX4600, and OCX Series switches.

Table 109: 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 110 on page 699](#) shows the default egress shared buffer allocation values in KB units.

Table 110: 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 111 on page 700](#) shows the default egress shared buffer allocation values for all switches as percentages.

Table 111: 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 112 on page 700](#) 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 112: Default Ingress and Egress Dedicated Buffer Pool Values (KB) per Switch (

Dedicated Buffer Type	QFX5210	QFX5200-48Y	QFX5110, QFX5200-32C	QFX5100, EX4600, OCX Series	QFX3500, 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 705](#).

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 113 on page 702](#) shows the default ingress shared buffer allocations:

Table 113: Default Ingress Shared Buffer Configuration

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

[Table 114 on page 702](#) shows the default egress shared buffer allocations:

Table 114: 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 115 on page 702](#)) and egress shared buffer settings (see [Table 116 on page 702](#)):

Table 115: 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 116: 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" on page 717](#) for an example that shows you how to configure the recommended buffer settings shown in [Table 115 on page 702](#) and [Table 116 on page 702](#).

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 117 on page 703](#)) and egress shared buffer settings (see [Table 118 on page 703](#)):

Table 117: 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 118: 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" on page 726](#) for an example that shows you how to configure the recommended buffer settings shown in [Table 115 on page 702](#) and [Table 116 on page 702](#).

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 119 on page 704](#)) and egress shared buffer settings (see [Table 120 on page 704](#)):

Table 119: 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 120: 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" on page 734](#) for an example that shows you how to configure the recommended buffer settings shown in [Table 119 on page 704](#) and [Table 120 on page 704](#).

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 121 on page 704](#)) and egress shared buffer settings (see [Table 122 on page 705](#)):

Table 121: 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 122: 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" on page 743](#) for an example that shows you how to configure the recommended buffer settings shown in [Table 121 on page 704](#) and [Table 122 on page 705](#).

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 multideestination), 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 | 717](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled | 726](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 734](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 743](#)

[Example: Configuring Queue Schedulers | 350](#)

[Configuring Global Ingress and Egress Shared Buffers | 708](#)

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, 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 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 | 717](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled | 726](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 734](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 743](#)

[Understanding CoS Buffer Configuration | 684](#)

Configuring Ingress and Egress Dedicated Buffers

SUMMARY

This topic describes how to adjust the size of the dedicated buffer, both globally and on a per-port basis.

IN THIS SECTION

- [Decreasing the Global Dedicated Buffer | 711](#)
- [Configuring and Applying Dedicated Buffer Profiles | 713](#)

The switch partitions its buffer into dedicated and shared buffers. As the name suggests, the dedicated buffer is exclusive to each port and only that port can use its dedicated buffer. The shared buffer is shared across all ports. When there is little traffic on many ports and few ports have bursty traffic, the dedicated buffers of the ports carrying little traffic are unused, and bursty traffic ports cannot use these unused buffers.

However, you can *decrease* the global dedicated buffer space from the default value, effectively increasing the global shared buffer space so that bursty traffic ports can use more of the buffer space according to their dynamic-threshold value.

You can also define a dedicated buffer profile to increase or decrease the dedicated buffer allocated to an individual port. This is particularly useful for decreasing dedicated buffer space on unused or down ports, thereby increasing dedicated buffer space available to active ports.

You can fine-tune the dedicated buffer settings, but do so with caution to prevent traffic loss due to buffer misconfiguration.

Decreasing the Global Dedicated Buffer

Set the global dedicated buffer pool of the device as a percent of its default. This, in effect, increases the shared buffer pool. The percentage can range from 15 to 100 percent of the default. The minimum of 15 percent is to ensure that each port receives some amount of minimal dedicated buffer to reduce the probability of all ports contending for the shared buffer.

NOTE: You can decrease the `dedicated-buffer` (thereby increasing the shared buffer pool), or decrease the `shared-buffer` (thereby increasing the dedicated buffer pool), but not both.

1. Check the current dedicated and shared buffer allocation.

For example:

```
user@host> show class-of-service shared-buffer
Ingress:
  Total Buffer      : 65536 KB
  Dedicated Buffer  : 7868 KB
  Shared Buffer     : 44420 KB
    Lossless       : 8884 KB
    Lossless Headroom : 4442 KB
    Lossy          : 31094 KB

Lossless Headroom Utilization:
Node Device      Total      Used      Free
0                4442 KB    0 KB     4442 KB
```

```
Egress:
  Total Buffer      : 65536 KB
  Dedicated Buffer   : 12739 KB
  Shared Buffer      : 44420 KB
  Lossless          : 8884 KB
  Lossy             : 31094 KB
```

2. Set the global egress dedicated buffer size as a percentage of its default, from 15 to 100 percent.

For example:

```
[edit class-of-service dedicated-buffer]
user@host# set egress percent 20
```

3. Set the global ingress dedicated buffer size as a percentage of its default, from 15 to 100 percent.

For example:

```
[edit class-of-service dedicated-buffer]
user@host# set ingress percent 25
```

4. Commit your changes.
5. Verify your configuration.

For example:

```
[edit class-of-service]
user@host# show
...
dedicated-buffer {
  ingress {
    percent 25;
  }
  egress {
    percent 20;
  }
}
...
```

6. Check the shared buffer to verify both the dedicated buffer and shared buffer have changed according to the configuration.

For example:

```

user@host> show class-of-service shared-buffer
Ingress:
  Total Buffer      : 65536 KB
  Dedicated Buffer : 1967 KB
  Shared Buffer    : 60513 KB
    Lossless       : 12102 KB
    Lossless Headroom : 6051 KB
    Lossy          : 42359 KB

  Lossless Headroom Utilization:
  Node Device      Total      Used      Free
  0                6051 KB    0 KB     6051 KB

Egress:
  Total Buffer      : 65536 KB
  Dedicated Buffer : 2547 KB
  Shared Buffer    : 60513 KB
    Lossless       : 12102 KB
    Lossy          : 42359 KB

```

Notice how the ingress and egress dedicated buffers are now less and the shared buffer has increased while the total buffer remains the same.

Configuring and Applying Dedicated Buffer Profiles

By default, the operating system calculates port level dedicated buffers internally. Therefore, even ports that are down or unused also get an equal amount of dedicated buffer space that is then unavailable to any traffic burst. The dedicated buffer profile provides the ability to increase or decrease the default dedicated buffer at a physical interface level.

With the dedicated buffer profile you can separately set the ingress and egress dedicated buffer size to be a number of cells, with each cell being 254 bytes. You can also set the ingress and egress dedicated buffer size to be none. Setting the dedicated buffer size to none is useful for unused or down ports. Note that when a port with no dedicated buffer becomes congested, the port directly consumes from the shared buffer pool.

Once you define a dedicated buffer profile, you can attach it directly to a physical interface.

Remaining dedicated buffers that you do not allocate to any ports by a dedicated buffer profile are equally shared among ports (based on speed) that don't have a dedicated buffer profile assigned to them.

NOTE: You cannot assign a `dedicated-buffer-profile` to aggregated Ethernet (ae-) interfaces. You can only assign a `dedicated-buffer-profile` to a physical interface.



CAUTION: If the buffer-size of all dedicated buffer profiles combined exceeds the total available dedicated buffer pool, the system logs a syslog error and does not implement the new configuration even though the commit succeeds.

1. Set the name and egress buffer size of the `dedicated-buffer-profile` .

For example:

```
[edit class-of-service]
user@host# set dedicated-buffer-profile dbp1 egress buffer-size 1000
```

2. Set the ingress buffer size of the dedicated buffer profile.

For example:

```
[edit class-of-service dedicated-buffer-profile dbp1]
user@host# set ingress buffer-size none
```

3. Commit your changes.
4. Apply the dedicated buffer profile to an interface.

For example:

```
[edit class-of-service]
user@host# set interfaces et-0/0/0 dedicated-buffer-profile dbp1
```

5. Verify your configuration.

For example:

```
[edit class-of-service]
user@host# show
...
dedicated-buffer-profile dbp1 {
    ingress {
        buffer-size {
            none;
        }
    }
}
```

```

    }
  }
  egress {
    buffer-size {
      1000;
    }
  }
}
...
interfaces {
  et-0/0/0 {
    dedicated-buffer-profile dbp1;
  }
}
...
```

6. Use show commands to verify the presence of the dedicated buffer profile.

For example:

```

user@host> show class-of-service dedicated-buffer-profile
Dedicated Buffer Profile: dbp1, Index: 1
Ingress Buffer Size: None
Egress Buffer Size: 1000
```

```

user@host> show class-of-service interface et-0/0/0
Physical interface: et-0/0/0, Index: 1004
Maximum usable queues: 10, Queues in use: 5
Exclude aggregate overhead bytes: disabled
Logical interface aggregate statistics: disabled
  Scheduler map: default, Index: 0
  Congestion-notification: Disabled
  Dedicated Buffer Profile: dbp1

  Logical interface: et-0/0/0.16386, Index: 1002
```

RELATED DOCUMENTATION

dedicated-buffer

dedicated-buffer-profile

show class-of-service dedicated-buffer-profile

Shared Buffer Pool Examples

IN THIS CHAPTER

- [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic | 717](#)
- [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled | 726](#)
- [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 734](#)
- [Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 743](#)

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

IN THIS SECTION

- [Requirements | 718](#)
- [Overview | 718](#)
- [Configuration | 720](#)
- [Verification | 723](#)

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 | 720](#)

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 123 on page 720 shows the configuration components for this example.

Table 123: 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% Percentage of ingress buffer space allocated to lossless traffic (lossless buffer partition): 5% Percentage of ingress buffer space allocated to lossless headroom traffic (lossless-headroom buffer partition): 0% 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% Percentage of egress buffer space allocated to lossless queues (lossless buffer partition): 5% Percentage of egress buffer space allocated to best-effort queues (lossy buffer partition): 75% Percentage of egress buffer space allocated to multicast traffic (multicast buffer partition): 20%

Configuration

IN THIS SECTION

- [CLI Quick Configuration | 721](#)
- [Configuring the Global Shared Buffer Pool for Networks with Mostly Best-Effort Unicast Traffic | 721](#)
- [Results | 722](#)

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 | 723](#)

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 | 726](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 734](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 743](#)

[Configuring Global Ingress and Egress Shared Buffers | 708](#)

[Understanding CoS Buffer Configuration | 684](#)

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 | 727](#)
- [Overview | 727](#)
- [Configuration | 729](#)
- [Verification | 732](#)

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 | 728](#)

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 124 on page 728 shows the configuration components for this example.

Table 124: Components of the Recommended Shared Buffer Configuration for Best-Effort Network Topologies with Links Enabled for Ethernet PAUSE

Component	Settings
Hardware	QFX3500 switch

Table 124: 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 | 729](#)
- [Configuring the Global Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links Enabled for Ethernet PAUSE | 730](#)
- [Results | 731](#)

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 | 732](#)

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 | 717](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 734](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 743](#)

[Configuring Global Ingress and Egress Shared Buffers | 708](#)

[Understanding CoS Buffer Configuration | 684](#)

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

IN THIS SECTION

- [Requirements | 735](#)
- [Overview | 735](#)
- [Configuration | 737](#)
- [Verification | 740](#)

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 | 737](#)

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 125 on page 737 shows the configuration components for this example.

Table 125: 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 | 738](#)
- [Configuring the Global Shared Buffer Pool for Networks with Mostly Multicast Traffic | 738](#)
- [Results | 739](#)

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 | 740](#)

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 | 717](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled | 726](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Lossless Traffic | 743](#)

[Configuring Global Ingress and Egress Shared Buffers | 708](#)

[Understanding CoS Buffer Configuration | 684](#)

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

IN THIS SECTION

- [Requirements | 744](#)
- [Overview | 744](#)
- [Configuration | 746](#)
- [Verification | 749](#)

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 | 745](#)

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 126 on page 746](#) shows the configuration components for this example.

Table 126: 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 | 747](#)
- [Configuring the Global Shared Buffer Pool for Networks with Mostly Lossless Traffic | 747](#)
- [Results | 748](#)

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:

```
rroot@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 | 749](#)

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 | 717](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Best-Effort Traffic on Links with Ethernet PAUSE Enabled | 726](#)

[Example: Recommended Configuration of the Shared Buffer Pool for Networks with Mostly Multicast Traffic | 734](#)

[Configuring Global Ingress and Egress Shared Buffers | 708](#)

[Understanding CoS Buffer Configuration | 684](#)



CoS on EVPN VXLANs

CoS Support on EVPN VXLANs | 753

CoS Support on EVPN VXLANs

IN THIS SECTION

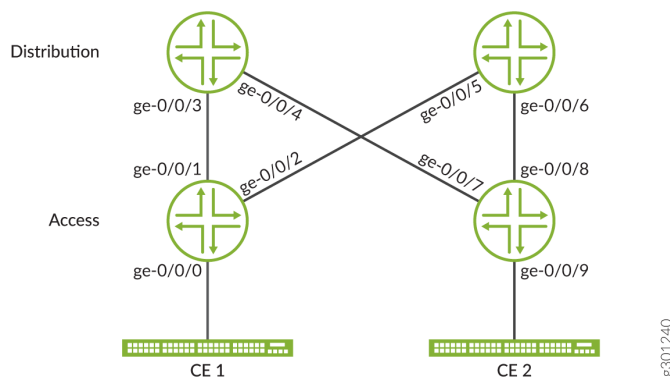
- Understanding CoS on VXLAN Interfaces | 753
- Configuring CoS on VXLAN Interfaces | 754
- Implementing CoS on VXLAN Interfaces (Junos OS Evolved) | 757
- CoS Limitations on VXLANs | 759

You can configure class of service (CoS) features on VXLAN interfaces. VXLAN traffic from different tenants traverses network boundaries over the same physical underlay network. To ensure fairness in the treatment of traffic for all tenants in the VXLAN, and to prioritize higher priority traffic, apply CoS features to the VXLAN interfaces.

Understanding CoS on VXLAN Interfaces

This section describes how classification and rewrite rules are applied to packets in a VXLAN instance. [Figure 28 on page 753](#) shows a simple VXLAN with two leaf nodes and one spine node.

Figure 28: Classifiers and Rewrite Rules on VXLANs



Refer to [Figure 28 on page 753](#) to understand the packet flow with DSCP/ToS fields in a VXLAN:

1. CE 1 sends a packet with Layer3 DSCP/ToS bit programmed to the Leaf 1 node.
2. Leaf 1 receives the original packet and appends the VXLAN header on top of the original packet. The outer VXLAN Layer3 header uses the original packet DSCP/Tos bit. You can create classifiers based on the original packet DSCP/802.1p bit. The ingress interface on the ingress leaf supports DSCP and 802.1p classifiers.
3. If rewrite is configured on Leaf 1, the inner header will have the DSCP/802.1p bit set by CE 1 and the outer header will have the rewrite bit. Only DSCP rewrite rules are supported, except on QFX10000 switches where 802.1p rewrite is also supported if the underlay is tagged.
4. The Spine node receives the VXLAN packet and can use ingress classification using these DSCP bits and forward the packet to the egress interface with the appropriate forwarding class.
5. The Spine egress interface can rewrite these bits using rewrite rules. These Spine rewrite rules only affects the outer Layer3 DSCP field. The inner/original packet still holds the DSCP/802.1p bit that was set by CE 1.
6. Leaf 2 receives the packet, processes the tunnel termination, and remove the outer VXLAN header.
7. Leaf 2 classification and rewrite functionality works on the inner header.
8. The original packet arrives on CE 2.

NOTE: On the leaf nodes, if the packet is multicast, you can use multi-destination classification to create appropriate multicast classification and rewrite rules.

Configuring CoS on VXLAN Interfaces

This section shows sample configurations of classifiers and rewrite rules for the leaf and spine nodes in VXLAN using [Figure 28 on page 753](#) as a reference. You can create schedulers as normal for the classifiers on each node.

Sample configuration of classifiers and rewrite rules on Leaf 1.

1. Create a classifier based on the *original*/DSCP/ToS bits:

```
[edit class-of-service classifiers]
user@leaf1#set dscp dscp_cf forwarding-class best-effort loss-priority low code-points 100000
user@leaf1#set dscp dscp_cf forwarding-class network-control loss-priority high code-points
110000
```

```

user@leaf1#set dscp dscp_cf forwarding-class expedited-forwarding loss-priority low code-
points 011010
user@leaf1#set dscp dscp_cf forwarding-class assured-forwarding loss-priority high code-
points 001010

```

2. Apply the classifier to the ingress interface:

```

[edit class-of-service interfaces]
user@leaf1#set ge-0/0/0 unit 0 classifiers dscp dscp_cf

```

3. Create a rewrite rule for the *outer* VXLAN DSCP/ToS bits:

```

[edit class-of-service rewrite-rules]
user@leaf1#set dscp dscp_rw forwarding-class best-effort loss-priority low code-points af22
user@leaf1#set dscp dscp_rw forwarding-class network-control loss-priority high code-points
af31
user@leaf1#set dscp dscp_rw forwarding-class expedited-forwarding loss-priority low code-
points af13
user@leaf1#set dscp dscp_rw forwarding-class assured-forwarding loss-priority high code-
points cs3

```

4. Apply the rewrite rule to the egress Leaf 1 interfaces:

```

[edit class-of-service interfaces]
user@leaf1#set ge-0/0/1 unit 0 rewrite-rules dscp dscp_rw
user@leaf1#set ge-0/0/2 unit 0 rewrite-rules dscp dscp_rw

```

Sample configuration of classifiers and rewrite rules on the Spine.

1. Create a classifier based on the outer VXLAN DSCP/ToS bits:

```

[edit class-of-service classifiers]
user@spine#set dscp dscp_cf forwarding-class best-effort loss-priority low code-points af22
user@spine#set dscp dscp_cf forwarding-class network-control loss-priority high code-points
af31
user@spine#set dscp dscp_cf forwarding-class expedited-forwarding loss-priority low code-
points af13
user@spine#set dscp dscp_cf forwarding-class assured-forwarding loss-priority high code-
points cs3

```


2. Apply the classifier to the ingress Spine interfaces:

```
[edit class-of-service interfaces]
user@spine#set ge-0/0/3 unit 0 classifiers dscp dscp_cf
user@spine#set ge-0/0/5 unit 0 classifiers dscp dscp_cf
```

3. Create a rewrite rule for the outer VXLAN DSCP/ToS bits:

```
[edit class-of-service rewrite-rules]
user@spine#set dscp dscp_rw forwarding-class best-effort loss-priority low code-points af22
user@spine#set dscp dscp_rw forwarding-class network-control loss-priority high code-points af31
user@spine#set dscp dscp_rw forwarding-class expedited-forwarding loss-priority low code-points af13
user@spine#set dscp dscp_rw forwarding-class assured-forwarding loss-priority high code-points cs3
```

4. Apply the rewrite rule to the egress Spine interfaces:

```
[edit class-of-service interfaces]
user@spine#set ge-0/0/4 unit 0 rewrite-rules dscp dscp_rw
user@spine#set ge-0/0/6 unit 0 rewrite-rules dscp dscp_rw
```

Sample configuration of classifiers and rewrite rules on Leaf 2.

1. Create a classifier based on the *original* DSCP/ToS bits, as the VXLAN header is removed at tunnel termination *before* forwarding classes are applied:

```
[edit class-of-service classifiers]
user@leaf2#set dscp dscp_cf forwarding-class best-effort loss-priority low code-points 100000
user@leaf2#set dscp dscp_cf forwarding-class network-control loss-priority high code-points 110000
user@leaf2#set dscp dscp_cf forwarding-class expedited-forwarding loss-priority low code-points 011010
user@leaf2#set dscp dscp_cf forwarding-class assured-forwarding loss-priority high code-points 001010
```

2. Apply the classifier to the ingress Leaf 2 interfaces:

```
[edit class-of-service interfaces]
user@leaf2#set ge-0/0/7 unit 0 classifiers dscp dscp_cf
user@leaf2#set ge-0/0/8 unit 0 classifiers dscp dscp_cf
```

3. Create a rewrite rule for the *original*/DSCP/ToS bits:

```
[edit class-of-service rewrite-rules]
user@leaf2#set dscp dscp_rw forwarding-class best-effort loss-priority low code-points 100000
user@leaf2#set dscp dscp_rw forwarding-class network-control loss-priority high code-points 110000
user@leaf2#set dscp dscp_rw forwarding-class expedited-forwarding loss-priority low code-points 011010
user@leaf2#set dscp dscp_rw forwarding-class assured-forwarding loss-priority high code-points 001010
```

4. Apply the rewrite rule to the egress Leaf 2 interface:

```
[edit class-of-service interfaces]
user@leaf2#set ge-0/0/9 unit 0 rewrite-rules dscp dscp_rw
```

To check the CoS configuration on one of the interfaces:

```
user@node#show class-of-service interface interface-name
```

To check the queue statistics on one of the interfaces:

```
user@node#show interfaces queue interface-name
```

Implementing CoS on VXLAN Interfaces (Junos OS Evolved)

CoS for EVPN VXLAN traffic is supported using a combination of classifiers, schedulers, and rewrite rules. This section describes how these components are implemented across different nodes on devices running Junos OS Evolved to apply CoS on the EVPN VXLAN traffic.

- **Classification at User Network Interface (UNI)/Ingress PE** — Traffic classification based on IEEE 802.1p and Differentiated Services code point (DSCP) are supported on the ingress PE where the EVPN VXLAN tunnel is initiated. BA and MF classifiers can be applied to Enterprise style (EP) or Service Provider (SP) style access interfaces.
- **Classification at Network Node Interface (NNI)/Egress PE** — Traffic classification based on IEEE 802.1p and Differentiated Services code point (DSCP) are supported on the egress PE where the EVPN VXLAN tunnel is terminated. BA classifiers can be applied to the underlying logical interface or unit. MF classifiers are not supported in tunnel terminations.
- **Rewrite at NNI** — After the encapsulation of the VXLAN tunnel, the rewrites on the outer/tunnel header are configured using the rewrite rules on the underlying logical interface or unit. Based on the configured rewrite rules, the VXLAN traffic is classified in the Spine/Network. DSCP rewrites on the outer/tunnel header of VXLAN packets is supported on the NNI interface.

Rewrite rules are supported in the following EVPN VXLAN scenarios:

- Intra-VNI L2 gateway — Rewrite rules are applied to both unicast and broadcast, unknown unicast, and multicast (BUM) traffic.
- Inter-VNI L3 gateway — Centrally-routed bridging (CRB) and edge-routed bridging (ERB).
- EVPN Type 5 routes.
- **Rewrite at UNI** — After the termination of the VXLAN tunnel, the rewrites on the inner headers are configured using rewrite rules on the Enterprise style (EP) or Service Provider (SP) style access interfaces. Based on the configured rewrite rules, the decapsulated packets are classified in the CE side network. The following rewrite rules are supported on the UNI interface for the decapsulated packets:
 - DSCP rewrites on the inner IPv4/IPv6 header
 - IEEE 802.1p rewrites on the inner Ethernet header (if tagged)

Rewrite rules are supported in the following EVPN VXLAN scenarios:

- Intra-VNI L2 gateway — Rewrite rules are applied to both unicast and broadcast, unknown unicast, and multicast (BUM) traffic.
- Inter-VNI L3 gateway — Centrally-routed bridging (CRB) and edge-routed bridging (ERB).
- EVPN Type 5 routes.
- **Scheduling** — Traffic prioritization and bandwidth reservation are achieved by using schedulers. The schedulers are associated with a forwarding class set via classifiers.

CoS Limitations on VXLANs

The following limitations apply to PTX routers:

- DSCP rewrite rules are not supported on Integrated Routing and Bridging (IRB) (L3 gateway scenarios).
- IEEE 802.1p rewrite rules are not supported on the NNI interface.
- Explicit congestion notification (ECN) rewrites are not supported on either UNI or NNI interfaces.
- Priority-based flow control (PFC) is not supported.
- No support for CoS classification and rewrite mechanism for IPv6 or IRB underlay.

The CoS functionality on EVPN VXLAN is the same as on QFX5K platforms. All VXLAN CoS features already supported on the QFX5120 are also supported on the QFX5130 and QFX5700 platforms.

The following limitations apply to the QFX5130 and QFX5700 platforms:

- HQoS is not supported due to hardware limitations.
- Classifier, rewrite and scheduler on IRB interface is not supported.
- DOT1P rewrite and classifier on the NNI port is not supported.
- DOT1P and DSCP rewrite on the UNI port is not supported.
- DSCP rewrite on the NNI port is supported with the following limitations:
 - DSCP rewrite takes effect only after you disable TOS copy (set `vxlan-disable-copy-tos-encap` at [edit forwarding-options] hierarchy level) on the VXLAN encapsulation node. When TOS copy is disabled, ECN bits are not copied from the inner to the outer header, so the packet outer header will have the defined rewrite DSCP value and an ECN value of 00.
 - DSCP rewrite rewrites both the outer and the inner header. So the inner header DSCP value cannot be preserved.
- PFC configuration will cause momentary traffic drops of up to 10ms.
- DSCP IPV6 classifiers and rewrites are not supported. Use DSCP classifier and rewrite instead.
- TOS copy feature does not work for Type-5 EVPN VXLANs.

The following limitation applies to QFX10000 platforms:

- Because IRB interfaces do not support dscp rewrite rules, you can apply rewrite rules on underlying L2 interfaces. 802.1p/dscp values in a VXLAN tunneled packet are written using underlying L2 interface rules.

7

PART

Configuration Statements and Operational Commands

[Junos CLI Reference Overview](#) | 761

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*