



Modular QoS Configuration Guide for Cisco 8000 Series Routers, IOS XR Release 7.5.x

First Published: 2021-11-01

Last Modified: 2023-09-01

Americas Headquarters

Cisco Systems, Inc.
170 West Tasman Drive
San Jose, CA 95134-1706
USA
<http://www.cisco.com>
Tel: 408 526-4000
800 553-NETS (6387)
Fax: 408 527-0883

THE SPECIFICATIONS AND INFORMATION REGARDING THE PRODUCTS IN THIS MANUAL ARE SUBJECT TO CHANGE WITHOUT NOTICE. ALL STATEMENTS, INFORMATION, AND RECOMMENDATIONS IN THIS MANUAL ARE BELIEVED TO BE ACCURATE BUT ARE PRESENTED WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED. USERS MUST TAKE FULL RESPONSIBILITY FOR THEIR APPLICATION OF ANY PRODUCTS.

THE SOFTWARE LICENSE AND LIMITED WARRANTY FOR THE ACCOMPANYING PRODUCT ARE SET FORTH IN THE INFORMATION PACKET THAT SHIPPED WITH THE PRODUCT AND ARE INCORPORATED HEREIN BY THIS REFERENCE. IF YOU ARE UNABLE TO LOCATE THE SOFTWARE LICENSE OR LIMITED WARRANTY, CONTACT YOUR CISCO REPRESENTATIVE FOR A COPY.

The Cisco implementation of TCP header compression is an adaptation of a program developed by the University of California, Berkeley (UCB) as part of UCB's public domain version of the UNIX operating system. All rights reserved. Copyright © 1981, Regents of the University of California.

NOTWITHSTANDING ANY OTHER WARRANTY HEREIN, ALL DOCUMENT FILES AND SOFTWARE OF THESE SUPPLIERS ARE PROVIDED "AS IS" WITH ALL FAULTS. CISCO AND THE ABOVE-NAMED SUPPLIERS DISCLAIM ALL WARRANTIES, EXPRESSED OR IMPLIED, INCLUDING, WITHOUT LIMITATION, THOSE OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NON-INFRINGEMENT OR ARISING FROM A COURSE OF DEALING, USAGE, OR TRADE PRACTICE.

IN NO EVENT SHALL CISCO OR ITS SUPPLIERS BE LIABLE FOR ANY INDIRECT, SPECIAL, CONSEQUENTIAL, OR INCIDENTAL DAMAGES, INCLUDING, WITHOUT LIMITATION, LOST PROFITS OR LOSS OR DAMAGE TO DATA ARISING OUT OF THE USE OR INABILITY TO USE THIS MANUAL, EVEN IF CISCO OR ITS SUPPLIERS HAVE BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

Any Internet Protocol (IP) addresses and phone numbers used in this document are not intended to be actual addresses and phone numbers. Any examples, command display output, network topology diagrams, and other figures included in the document are shown for illustrative purposes only. Any use of actual IP addresses or phone numbers in illustrative content is unintentional and coincidental.

All printed copies and duplicate soft copies of this document are considered uncontrolled. See the current online version for the latest version.

Cisco has more than 200 offices worldwide. Addresses and phone numbers are listed on the Cisco website at www.cisco.com/go/offices.

Cisco and the Cisco logo are trademarks or registered trademarks of Cisco and/or its affiliates in the U.S. and other countries. To view a list of Cisco trademarks, go to this URL: <https://www.cisco.com/c/en/us/about/legal/trademarks.html>. Third-party trademarks mentioned are the property of their respective owners. The use of the word partner does not imply a partnership relationship between Cisco and any other company. (1721R)

© 2021–2022 Cisco Systems, Inc. All rights reserved.



CONTENTS

PREFACE

Preface ix

Changes to This Document **ix**

Communications, Services, and Additional Information **ix**

CHAPTER 1

New and Changed QoS Features 1

New and Changed QoS Features **1**

CHAPTER 2

Traffic Management Overview 3

Scope **3**

Traditional Traffic Management **3**

Traffic Management on Your Router **3**

Limitations of the VoQ Model **4**

QoS Policy Inheritance **5**

Cisco Modular QoS CLI to Deploy QoS **6**

Important Points about MQC Egress Queuing Policy **6**

CHAPTER 3

Classify Packets to Identify Specific Traffic 9

Classify Packets to Identify Specific Traffic **9**

Packet Classification Overview **9**

Specification of the CoS for a Packet with IP Precedence **10**

IP Precedence Bits Used to Classify Packets **10**

IP Precedence Value Settings **10**

IP Precedence Compared to IP DSCP Marking **11**

Packet Classification on Your Router **11**

Classify and Remark Layer 3 Header on Layer 2 Interfaces **11**

Traffic Class Elements **12**

- Default Traffic Class 13
- Create a Traffic Class 14
- Traffic Policy Elements 15
- Create a Traffic Policy 16
- Attach a Traffic Policy to an Interface 17
- QoS Policy Propagation via BGP (QPPB) 20
 - Configuration Workflow 22
 - Egress Interface Configuration 24

CHAPTER 4

Mark Packets to Change Priority Settings 27

- Packet Marking Overview 27
 - Default Marking 27
 - QoS Behavior for Generic Routing Encapsulation (GRE) Tunnels 28
 - Packet Marking 28
 - QoS Behavior for Generic Routing Encapsulation (GRE) Tunnels 29
- Class-based Unconditional Packet Marking Feature and Benefits 29
- Configure Class-based Unconditional Packet Marking 30
- Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic 32
 - Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic: Benefits 33
 - Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic: Guidelines and Limitations 33
 - Configure IP DSCP and MPLS EXP Values for IP to MPLS Traffic 33
- Class-based Unconditional Packet Marking: Examples 35
 - IP Precedence Marking Configuration: Example 35
 - IP DSCP Marking Configuration: Example 35
 - QoS Group Marking Configuration: Example 36
 - CoS Marking Configuration: Example 36
 - MPLS Experimental Bit Imposition Marking Configuration: Example 37
 - MPLS Experimental Topmost Marking Configuration: Example 37
 - Queuing and Marking Policies: Examples 37
- IP Precedence Compared to IP DSCP Marking 40
 - Configure DSCP CS7 (Precedence 7) 40
- In-Place Policy Modification 41
 - Recommendations for Using In-Place Policy Modification 41

CHAPTER 5**Congestion Avoidance 43**

Congestion Avoidance 43

Queuing Modes 43

Main Interface Queuing Policy 44

Sub-Interface Queuing Policy 44

Congestion Avoidance in VOQ 44

Sharing of VOQ Statistics Counters 45

Configuring Sharing of VOQ Statistics Counters 45

Dual Queue Limit 46

Restrictions 47

Equitable Traffic Flow Using Fair VOQ 48

Fair VOQ: Why 48

Fair VOQ: How 49

Fair VOQ Modes and Sharing of Counters 50

Fair VOQs and Slice (or Normal) VOQs: Key Differences 51

Guidelines and Limitations 51

Configure Fair VOQ 53

Modular QoS Congestion Avoidance 54

Tail Drop and the FIFO Queue 54

Configure Tail Drop 54

Random Early Detection and TCP 56

Configure Random Early Detection 56

Explicit Congestion Notification 58

CHAPTER 6**Configure Priority Flow Control 63**

Priority Flow Control Overview 63

buffer-internal mode 65

Restrictions and Guidelines 65

buffer-extended mode 65

Important Considerations 66

Hardware Support for Priority Flow Control 67

Configure Priority Flow Control 67

View Priority Flow Control TX Frames Per Traffic Class 72

- Guidelines and Limitations for Viewing Priority Flow Control TX Frames Per Traffic Class 73
- Configurable ECN Threshold and Maximum Marking Probability Values 74
 - ECN Threshold and Maximum Marking Probability Values 74
 - Benefits of Configurable ECN Threshold and Maximum Marking Probability Values 75
 - ECN Threshold and Maximum Marking Probability Values: FAQs 75
 - Guidelines and Limitations 76
 - Configure ECN Threshold and Maximum Marking Probability Values 77
- Priority Flow Control Watchdog Overview 79
 - Configure a Priority Flow Control Watchdog Interval 80
- Global Statistics Counters for Priority Flow Control and Priority Flow Control Watchdog 82
 - Global Statistics Counters for Priority Flow Control and Priority Flow Control Watchdog: Benefits 83
 - Global Statistics Counters for Priority Flow Control and Priority Flow Control Watchdog: Guidelines and Limitations 83
 - View Global Statistics for Priority Flow Control and Priority Flow Control Watchdog 83
- PFC Watchdog Statistics to Monitor Traffic Drops 85
 - PFC Watchdog Statistics for Two Counters 86
 - PFC Watchdog Statistics Benefits 86
 - View PFC Watchdog Statistics 86
- Detect High Bandwidth Memory Congestion 87
 - Benefits of Detecting High Bandwidth Memory Congestion 88
 - High Bandwidth Memory Congestion Detection: Guidelines and Limitations 88
 - Configure High Bandwidth Memory Congestion Detection 88
 - Available Shared Memory System and High Bandwidth Memory 93
 - Global Pause Frames for High Bandwidth Memory Congestion 94
 - Global Pause Frames for High Bandwidth Memory Congestion: Guidelines and Limitations 95
 - Configure Global Pause Frames for High Bandwidth Memory Congestion 95

CHAPTER 7

Congestion Management 97

- Congestion Management Overview 97
- Low-Latency Queueing with Strict Priority Queueing 97
 - Configure Low Latency Queueing with Strict Priority Queueing 97
- Traffic Shaping 100
 - Configure Traffic Shaping 100
- Traffic Policing 102

Committed Bursts and Excess Bursts	102
Single-Rate Policer	103
Two-Rate Policer	105
Committed Bursts	107
Excess Bursts	108
Two-Rate Policer Details	109
References for Modular QoS Management	109

CHAPTER 8**Configure Modular QoS on Link Bundles 111**

QoS on Link Bundles	111
Load Balancing	111
Configure QoS on Link Bundles	112



Preface

This preface contains these sections:

- [Changes to This Document, on page ix](#)
- [Communications, Services, and Additional Information, on page ix](#)

Changes to This Document

This table lists the technical changes made to this document since it was first published.

Table 1: Changes to this Document

Date	Change Summary
March 2023	Republished for Release 7.5.4
September 2022	Republished for Release 7.5.3
April 2022	Republished for Release 7.5.2
November 2021	Initial release of this document.

Communications, Services, and Additional Information

- To receive timely, relevant information from Cisco, sign up at [Cisco Profile Manager](#).
- To get the business impact you're looking for with the technologies that matter, visit [Cisco Services](#).
- To submit a service request, visit [Cisco Support](#).
- To discover and browse secure, validated enterprise-class apps, products, solutions and services, visit [Cisco Marketplace](#).
- To obtain general networking, training, and certification titles, visit [Cisco Press](#).
- To find warranty information for a specific product or product family, access [Cisco Warranty Finder](#).

Cisco Bug Search Tool

[Cisco Bug Search Tool](#) (BST) is a web-based tool that acts as a gateway to the Cisco bug tracking system that maintains a comprehensive list of defects and vulnerabilities in Cisco products and software. BST provides you with detailed defect information about your products and software.



CHAPTER 1

New and Changed QoS Features

- [New and Changed QoS Features, on page 1](#)

New and Changed QoS Features

Table 2: QoS Features Added or Modified in IOS XR Release 7.5.x

Feature	Description	Changed in Release	Where Documented
PFC Watchdog Statistics to Monitor Traffic Drops	This feature was introduced.	Release 7.5.4	PFC Watchdog Statistics to Monitor Traffic Drops
Global Pause Frames for High Bandwidth Memory Congestion	This feature was introduced.	Release 7.5.4	Global Pause Frames for High Bandwidth Memory Congestion
View Priority Flow Control TX Frames Per Traffic Class	This feature was introduced.	Release 7.5.4	View Priority Flow Control TX Frames Per Traffic Class
Configurable ECN Threshold and Maximum Marking Probability Values	This feature was introduced.	Release 7.5.4	Configurable ECN Threshold and Maximum Marking Probability Values, on page 74
Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic	This feature was introduced.	Release 7.5.4	Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic
Priority Flow Control on Cisco 8808 and Cisco 8812 Modular Chassis Line Cards	Priority Flow Control is now supported on additional Cisco 8808 and Cisco 8812 Modular Chassis Line Cards .	Release 7.5.3	Priority Flow Control Overview, on page 63

Feature	Description	Changed in Release	Where Documented
Detect High Bandwidth Memory Congestion	We provide detailed insights into congestion on the High Bandwidth Memory (HBM) write channel, such as the devices on which congestion has occurred, the time stamps, and when the device returned to its normal state.	Release 7.5.3	Detect High Bandwidth Memory Congestion, on page 87
QoS Policy Propagation via BGP	You now have the ability to install a BGP route in the routing table with a QoS Group so that IP packets that match the route receive the QoS policies associated with the QoS group.	Release 7.5.2	QoS Policy Propagation via BGP (QPPB), on page 20



CHAPTER 2

Traffic Management Overview

- [Scope, on page 3](#)
- [Traditional Traffic Management, on page 3](#)
- [Traffic Management on Your Router, on page 3](#)
- [Limitations of the VoQ Model, on page 4](#)
- [QoS Policy Inheritance, on page 5](#)
- [Cisco Modular QoS CLI to Deploy QoS, on page 6](#)

Scope

Read this configuration guide to understand the overall architecture that powers the Cisco Quality of Service (QoS) technology, and also how to use its features to configure and manage the traffic bandwidth and packet loss parameters on your network.

Traditional Traffic Management

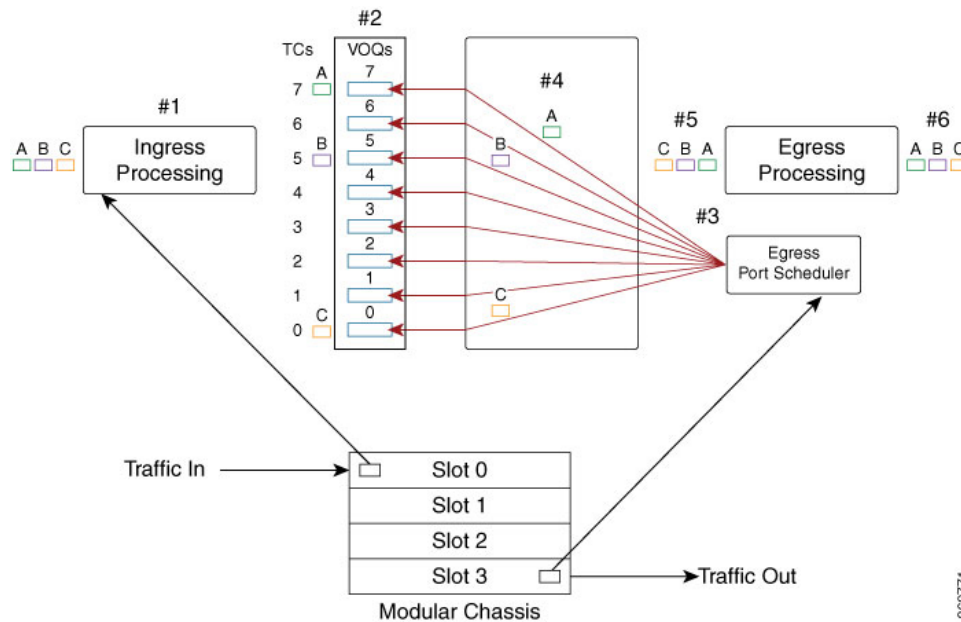
In traditional methods of traffic management, traffic packets are sent to the egress output queues without taking into consideration the egress interface availability to transmit.

Therein lies the problem as well. In case of a traffic congestion, the traffic packets may get dropped at the egress port. Which means that the network resources spent getting the packets from the ingress input queue across the switch fabric to the output queues at egress are wasted. That is not all—input queues buffer traffic meant for different egress ports, so congestion on one egress port could affect traffic on another port, an event referred to as head-of-line-blocking.

Traffic Management on Your Router

Your router's Network Processing Unit (NPU) uses a coupled ingress-egress virtual output queuing (VoQ)-based forwarding architecture to manage traffic.

Figure 1: Traffic flow from ingress port on slice 0 to an egress port on slot 3



Here, each ingress traffic class has a one-to-one VoQ mapping from each ingress slice (pipeline) to each egress port. Which means that every egress interface (#5 in the figure) has earmarked buffer space on every ingress pipeline (#1 in the figure) for each of its VoQs.

Here is how the story of packet travel in times of congestion on your router system unravels:

#1: Packets A (colored green), B (colored pink), and C (colored brown) are at the ingress interface. This is where packet marking, classification, and policing takes place. (For details, see [Mark Packets to Change Priority Settings, on page 27](#), [Classify Packets to Identify Specific Traffic, on page 9](#), and [Congestion Management, on page 97](#).)

#2: These packets are stored in separate buffer storage spaces in dedicated VoQs. This is where queuing, VoQ transmit, and drop packet and byte counters come into play. (For details, see [Congestion Avoidance, on page 43](#).)

#3: Depending on the bandwidth available on the egress interface, these packets are subjected to egress scheduling, where egress credit and transmit schedulers are configured. In other words, the packets and the sequence in which they will now proceed towards the egress interface is determined here. This is where the fabric bandwidth is taken into consideration for egress scheduling.

#4: The packets are switched through the fabric.

#5: In the final phase, the egress marking and classification takes place, and the congestion is managed in a way that at this stage there is no packet dropped, and all the packets are transmitted to the next hop.

Limitations of the VoQ Model

While the VoQ model of traffic management offers distinct advantages (reducing memory bandwidth requirements, providing end-to-end QoS flow), it has this limitation:

The total egress queue scale is lower because each egress queue must be replicated as an ingress VoQ on each slice of each NPU/ASIC in the system. This means that with the addition of 1 NPU with 20 interfaces, the

number of VoQs used on each and every NPU in the system will increase by 20×8 (queue/interface) = 160. There is also an increase in the number of credit connectors from each scheduler for each egress port on pre-existing NPUs to each slice in the newly inserted NPU.

QoS Policy Inheritance

Table 3: Feature History Table

Feature Name	Release Information	Feature Description
QoS Policy Inheritance	Release 7.3.15	<p>To create QoS policies for subinterfaces, you had to apply the policy on each subinterface manually. From this release, all you do is create and apply a single QoS policy on the main interface and the subinterfaces automatically inherit the policy.</p> <p>The inheritance model provides an easily maintainable method for applying policies, enabling you to create targeted policies for a group of interfaces and their subinterfaces. This model saves your time and resources while creating QoS policies.</p>

- **What is this functionality all about?**—As the name suggests, the functionality is based on an inheritance model, where you create and apply a QoS policy to a main interface. The subinterfaces that are attached to the main interface automatically inherit the policy.

The inheritance model applies to all QoS operations including:

- Classification
- Marking
- Policing
- Shaping

- **How does the inheritance model help?**—Previously, if you had, for example, eight subinterfaces, you created and applied policies to each of those subinterfaces separately. With the inheritance model, you save on time and resources, with just one policy automatically applied across a main interface and its subinterfaces.
- **Do I need to do anything more to enable the inheritance model?**—No, you don't. The inheritance model is the default option.
- **What if I want to override the inheritance option?**—Technically, you can't override this option. However, you can remove the policy from main interface and add policies to the subinterfaces except to the ones where you don't want the policy to be inherited.

- **What about policy-map statistics?**—There's no change to this behavior. Running the [show policy-map interface](#) command displays the cumulative statistics for an interface, and these numbers include the subinterfaces as well.
- **Any limitations I need to be aware of?**—There's no support for ECN marking and egress marking policy on the same interface and subinterface combination. However, the QoS policy inheritance functionality accepts these multiple policies causing the ECN markings to fail. To prevent such failures:
 - Don't configure an egress marking policy on a subinterface and apply an ECN-enabled policy on the main interface.
 - Don't apply an ECN policy on a subinterface and configure an egress marking policy on the main interface.

Cisco Modular QoS CLI to Deploy QoS

Cisco Modular QoS CLI (MQC) framework is the Cisco IOS QoS user language that enables:

- a standard Command Line Interface (CLI) and semantics for QoS features.
- simple and accurate configurations.
- QoS provisioning within the context of an extensible language.

For your router, in the egress direction, two types of MQC policies are supported: queuing and marking. You use the queuing policy to configure credit scheduling hierarchy, rates, priority, buffering, and congestion avoidance. You use the marking policy to classify and mark packets that have been scheduled for transmission. Even when a queuing policy is not applied, there is an implicit queuing policy with TC7 - P1, TC6 - P2, TC5 - TC0 (6 x Pn), so packets marked with TC7 and control inject packets are always prioritized over other packets. In the ingress, only one policy is supported for classification and marking.

You can apply queuing and marking policy independent of each other or together in the egress direction. If you apply both policies together, the queuing policy actions are provisioned first, followed by marking policy actions.

Important Points about MQC Egress Queuing Policy

These are important points that you must know about the MQC egress queuing policy:

- The MQC queuing policy consists of a set of class maps, which are added to a policy map. You control queuing and scheduling parameters for that traffic class by applying actions to the policy.
- **class-default** matches **traffic-class 0**, and no other class can match **traffic-class 0**. When **set traffic-class** is used in the ingress policy-map, then it should match the **traffic-class** in the egress queuing policy-map to have a predictable behaviour.
- Each unique combination of traffic classes that match **class-default** require a separate traffic class (TC) profile. The number of TC profiles are limited to 8 for main interfaces and 8 for sub-interfaces.
- You cannot configure multiple traffic classes with the same priority level.
- Each priority level, when configured, must be configured to the class that matches the corresponding TC as shown in the following table.

Priority Level	Traffic Class
P1	7
P2	6
P3	5
P4	4
P5	3
P6	2
P7	1

- If all the priority levels configured in a policy-map are sorted, they must be contiguous. In other words, you cannot skip a priority level. For example, P1 P2 P4 (skipping P3), is not allowed.
- From IOS XR Release 7.3.1 onwards, you can create a single set of contiguous priority TCs. Ensure that you assign priority levels that increase for every TC or remain the same, but don't decrease. Also, make sure that you assign priority level 1 for traffic class 7. You need not configure unused traffic classes, so you can create only those many TCs that you require on the egress policy-map.
- MQC supports up to two levels (parent, child) of queuing policy. The parent level aggregates all the traffic classes and whereas the child level differentiates traffic classes using MQC classes.
- Only these actions are supported in the queuing policy:
 - priority
 - shape
 - bandwidth remaining ratio
 - queue-limit
 - Random Early Detection (RED)
 - Priority flow control
- You can have only one match traffic-class value in the class map.
- You cannot apply a queuing policy to a main interface and its sub-interfaces.



CHAPTER 3

Classify Packets to Identify Specific Traffic

- [Classify Packets to Identify Specific Traffic, on page 9](#)
- [Packet Classification Overview, on page 9](#)
- [Packet Classification on Your Router, on page 11](#)
- [Traffic Class Elements, on page 12](#)
- [Default Traffic Class, on page 13](#)
- [Create a Traffic Class, on page 14](#)
- [Traffic Policy Elements, on page 15](#)
- [Create a Traffic Policy, on page 16](#)
- [Attach a Traffic Policy to an Interface, on page 17](#)
- [QoS Policy Propagation via BGP \(QPPB\), on page 20](#)

Classify Packets to Identify Specific Traffic

Read this section to get an overview of packet classification and the different packet classification types for your router.

Packet Classification Overview

Packet classification involves categorizing a packet within a specific group (or class) and assigning it a traffic descriptor to make it accessible for QoS handling on the network. The traffic descriptor contains information about the forwarding treatment (quality of service) that the packet should receive. Using packet classification, you can partition network traffic into multiple priority levels or classes of service.

When traffic descriptors are used to classify traffic, the source agrees to adhere to the contracted terms and the network promises a quality of service. This is where traffic policers and traffic shapers come into the picture. Traffic policers and traffic shapers use the traffic descriptor of a packet—that is, its classification—to ensure adherence to the contract.

The Modular Quality of Service (QoS) command-line interface (MQC) is used to define the traffic flows that must be classified, where each traffic flow is called a class of service, or class. Later, a traffic policy is created and applied to a class. All traffic not identified by defined classes fall into the category of a default class.

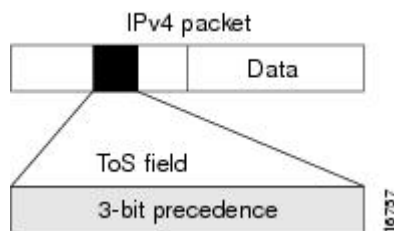


Note From Cisco IOS XR Release 7.2.12 onwards, you can classify packets on Layer 2 transport interfaces using Layer 3 header values. However, this feature applies only to the main interface (physical and bundle interfaces), and not on the sub-interfaces.

Specification of the CoS for a Packet with IP Precedence

Use of IP precedence allows you to specify the CoS for a packet. You can create differentiated service by setting precedence levels on incoming traffic and using them in combination with the QoS queuing features. So that, each subsequent network element can provide service based on the determined policy. IP precedence is usually deployed as close to the edge of the network or administrative domain as possible. This allows the rest of the core or backbone to implement QoS based on precedence.

Figure 2: IPv4 Packet Type of Service Field



You can use the three precedence bits in the type-of-service (ToS) field of the IPv4 header for this purpose. Using the ToS bits, you can define up to eight classes of service. Other features configured throughout the network can then use these bits to determine how to treat the packet in regard to the ToS to grant it. These other QoS features can assign appropriate traffic-handling policies, including congestion management strategy and bandwidth allocation. For example, queuing features such as LLQ can use the IP precedence setting of the packet to prioritize traffic.

IP Precedence Bits Used to Classify Packets

Use the three IP precedence bits in the ToS field of the IP header to specify the CoS assignment for each packet. You can partition traffic into a maximum of eight classes and then use policy maps to define network policies in terms of congestion handling and bandwidth allocation for each class.

Each precedence corresponds to a name. IP precedence bit settings 6 and 7 are reserved for network control information, such as routing updates. These names are defined in RFC 791.

IP Precedence Value Settings

By default, the routers leave the IP precedence value untouched. This preserves the precedence value set in the header and allows all internal network devices to provide service based on the IP precedence setting. This policy follows the standard approach stipulating that network traffic should be sorted into various types of service at the edge of the network and that those types of service should be implemented in the core of the network. Routers in the core of the network can then use the precedence bits to determine the order of transmission, the likelihood of packet drop, and so on.

Because traffic coming into your network can have the precedence set by outside devices, we recommend that you reset the precedence for all traffic entering your network. By controlling IP precedence settings, you

prohibit users that have already set the IP precedence from acquiring better service for their traffic simply by setting a high precedence for all of their packets.

The class-based unconditional packet marking and LLQ features can use the IP precedence bits.

IP Precedence Compared to IP DSCP Marking

If you need to mark packets in your network and all your devices support IP DSCP marking, use the IP DSCP marking to mark your packets because the IP DSCP markings provide more unconditional packet marking options. If marking by IP DSCP is undesirable, however, or if you are unsure if the devices in your network support IP DSCP values, use the IP precedence value to mark your packets. The IP precedence value is likely to be supported by all devices in the network.

You can set up to 8 different IP precedence markings and 64 different IP DSCP markings.

Packet Classification on Your Router

On your router, there are two types of packet classification systems:

- In the ingress direction, QoS map and Ternary Content Addressable Memory (TCAM).



Note TCAM is not supported on fixed-configuration routers (where router interfaces are built in). It is supported only on modular routers (that have multiple slots that allow you to change the interfaces on the router).

- In the egress direction, QoS map.

When a policy is matching only on Differentiated Services Code Point (DSCP) or precedence value (also called DSCP or Precedence-based classification), the system selects map-based classification system; else, it selects TCAM.

The TCAM is an extension of the Content Addressable Memory (CAM) table concept. A CAM table takes in an index or key value (usually a MAC address) and looks up the resulting value (usually a switch port or VLAN ID). Table lookup is fast and always based on an exact key match consisting of two input values: 0 and 1 bits.

The QoS map is a table-based classification system for traffic packets.

Classify and Remark Layer 3 Header on Layer 2 Interfaces

When you need to mark packets for Layer 2 interface traffic that flows across bridge domains and bridge virtual interfaces (BVI), you can create a mixed QoS policy. This policy has both map-based and TCAM-based classification class-maps. The mixed policy ensures that both bridged (Layer 2) and Bridge Virtual Interface (BVI, or Layer 3) traffic flows are classified and remarked.

Guidelines

- A class-map with TCAM classification may not match bridged traffic. TCAM entries match only routed traffic while map entries match both bridged and BVI traffic.

- A class-map with map-based classification matches both bridged and BVI traffic.

Example

```

ipv4 access-list acl_v4
10 permit ipv4 host 100.1.1.2 any
20 permit ipv4 host 100.1.100.2 any
ipv6 access-list acl_v6
10 permit tcp host 50:1:1::2 any
20 permit tcp any host 50:1:200::2
class-map match-any c_match_acl
match access-group ipv4 acl_v4 ! This entry does not match bridged traffic
match access-group ipv6 acl_v6 ! This entry does not match bridged traffic
match dscp af11 This entry matches bridged and BVI traffic
class-map match-all c_match_all
match protocol udp ! This entry does not match bridged traffic
match prec 7
class-map match-any c_match_protocol
match protocol tcp ! This entry, and hence this class does not match bridged traffic
class-map match-any c_match_ef
match dscp ef ! This entry/class matches bridged and BVI traffic
class-map match-any c_qosgroup_1 This class matches bridged and BVI traffic
!
match qos-group 1
policy-map p_ingress
class c_match_acl
set traffic-class 1
set qos-group 1
!
class c_match_all
set traffic-class 2
set qos-group 2
!
class c_match_ef
set traffic-class 3
set qos-group 3
!
class c_match_protocol
set traffic-class 4
set qos-group 4
policy-map p_egress
class c_qosgroup_1
set dscp af23
interface FourHundredGigE0/0/0/0
l2transport
service-policy input p_ingress
service-policy output p_egress
!
!
interface FourHundredGigE0/0/0/1
ipv4 address 200.1.2.1 255.255.255.0
ipv6 address 2001:2:2::1/64
service-policy input p_ingress
service-policy output p_egress

```

Traffic Class Elements

The purpose of a traffic class is to classify traffic on your router. Use the **class-map** command to define a traffic class.

A traffic class contains three major elements:

- A name
- A series of **match** commands - to specify various criteria for classifying packets.
- An instruction on how to evaluate these **match** commands (if more than one **match** command exists in the traffic class)

Packets are checked to determine whether they match the criteria specified in the **match** commands. If a packet matches the specified criteria, that packet is considered a member of the class and is forwarded according to the QoS specifications set in the traffic policy. Packets that fail to meet any of the matching criteria are classified as members of the default traffic class.

This table shows the details of match types supported on the router.

Match Type Supported	Min, Max	Max Entries	Support for Match NOT	Support for Ranges	Direction Supported on Interfaces
IPv4 DSCP IPv6 DSCP	(0,63)	64	Yes	Yes	Ingress
DSCP					Egress
IPv4 Precedence IPv6 Precedence	(0,7)	8	Yes	No	Ingress
Precedence					Egress
MPLS Experimental Topmost	(0,7)	8	Yes	No	Ingress Egress
Access-group	Not applicable	8	No	Not applicable	Ingress
QoS-group	(1,7)	7	No	No	Egress
Protocol	(0, 255)	1	Yes	Not applicable	Ingress

Default Traffic Class

Unclassified traffic (traffic that does not meet the match criteria specified in the traffic classes) is treated as belonging to the default traffic class.

If the user does not configure a default class, packets are still treated as members of the default class. However, by default, the default class has no enabled features. Therefore, packets belonging to a default class with no configured features have no QoS functionality.

For egress classification, match on **qos-group** (1-7) is supported. Match **qos-group 0** cannot be configured. The class-default in the egress policy maps to **qos-group 0**.

This example shows how to configure a traffic policy for the default class:

```
configure
policy-map ingress_policy1
class class-default
  police rate percent 30
!
```

Create a Traffic Class

To create a traffic class containing match criteria, use the **class-map** command to specify the traffic class name, and then use the **match** commands in class-map configuration mode, as needed.

Guidelines

- Users can provide multiple values for a match type in a single line of configuration; that is, if the first value does not meet the match criteria, then the next value indicated in the match statement is considered for classification.
- Use the **not** keyword with the **match** command to perform a match based on the values of a field that are not specified.
- All **match** commands specified in this configuration task are considered optional, but you must configure at least one match criterion for a class.
- If you specify **match-any**, one of the match criteria must be met for traffic entering the traffic class to be classified as part of the traffic class. This is the default. If you specify **match-all**, the traffic must match all the match criteria.
- For the **match access-group** command, QoS classification based on the packet length or TTL (time to live) field in the IPv4 and IPv6 headers is not supported.
- For the **match access-group** command, when an ACL list is used within a class-map, the deny action of the ACL is ignored and the traffic is classified based on the specified ACL match parameters.
- The **match qos-group**, **traffic-class**, **DSCP/Prec**, and **MPLS EXP** are supported only in egress direction, and these are the only match criteria supported in egress direction.
- The egress default class implicitly matches **qos-group 0**.
- Multicast takes a system path that is different than unicast on router, and they meet later on the egress in a multicast-to-unicast ratio of 20:80 on a per interface basis. This ratio is maintained on the same priority level as that of the traffic.
- Egress QoS for multicast traffic treats traffic classes 0-5 as low-priority and traffic classes 6-7 as high priority. Currently, this is not user-configurable.
- Egress shaping does not take effect for multicast traffic in the high priority (HP) traffic classes. It only applies to unicast traffic.
- If you set a traffic class at the ingress policy and do not have a matching class at egress for the corresponding traffic class value, then the traffic at ingress with this class will not be accounted for in the default class at the egress policy map.

- Only traffic class 0 falls in the default class. A non-zero traffic class assigned on ingress but with no assigned egress queue, falls neither in the default class nor any other class.

Configuration Example

You have to accomplish the following to complete the traffic class configuration:

1. Creating a class map
2. Specifying the match criteria for classifying the packet as a member of that particular class
(For a list of supported match types, see [Traffic Class Elements, on page 12.](#))

```
Router# configure
Router(config)# class-map match-any qos-1
Router(config-cmap)# match qos-group 1
Router(config-cmap)# end-class-map
Router(config-cmap)# commit
```

Use this command to verify the class-map configuration:

```
Router#show class-map qos-1
1) ClassMap: qos-1      Type: qos
   Referenced by 2 Policymaps
```

Also see, [Attach a Traffic Policy to an Interface, on page 17.](#)

Related Topics

- [Traffic Class Elements, on page 12](#)
- [Traffic Policy Elements, on page 15](#)

Traffic Policy Elements

A traffic policy contains three elements:

- Name
- Traffic class
- QoS policies

After choosing the traffic class that is used to classify traffic to the traffic policy, the user can enter the QoS features to be applied to the classified traffic.

The MQC does not necessarily require that the users associate only one traffic class to one traffic policy.

The order in which classes are configured in a policy map is important. The match rules of the classes are programmed into the TCAM in the order in which the classes are specified in a policy map. Therefore, if a packet can possibly match multiple classes, only the first matching class is returned and the corresponding policy is applied.

The router supports 8 classes per policy-map in the ingress direction and 8 classes per policy-map in the egress direction.

This table shows the supported class-actions on the router.

Supported Action Types	Direction supported on Interfaces
bandwidth-remaining	egress
mark	See Packet Marking, on page 28
police	ingress
priority	egress (level 1 to level 7)
queue-limit	egress
shape	egress
red	egress

RED supports the **discard-class** option; the only values to be passed to the discard-class being 0 and 1.

Create a Traffic Policy

The purpose of a traffic policy is to configure the QoS features that should be associated with the traffic that has been classified in a user-specified traffic class or classes.

To configure a traffic class, see [Create a Traffic Class, on page 14](#).

After you define a traffic policy with the **policy-map** command, you can attach it to one or more interfaces to specify the traffic policy for those interfaces by using the **service-policy** command in interface configuration mode. With dual policy support, you can have two traffic policies, one marking and one queuing attached at the output. See, [Attach a Traffic Policy to an Interface, on page 17](#).

Configuration Example

You have to accomplish the following to complete the traffic policy configuration:

1. Creating a policy map that can be attached to one or more interfaces to specify a service policy
2. Associating the traffic class with the traffic policy
3. Specifying the class-action(s) (see [Traffic Policy Elements, on page 15](#))

```
Router# configure
Router(config)# policy-map test-shape-1
Router(config-pmap)# class qos-1

/* Configure class-action ('shape' in this example).
Repeat as required, to specify other class-actions */
Router(config-pmap-c)# shape average percent 40
Router(config-pmap-c)# exit

/* Repeat class configuration as required, to specify other classes */
```

```
Router(config-pmap)# end-policy-map
Router(config)# commit
```

Related Topics

- [Traffic Policy Elements, on page 15](#)
- [Traffic Class Elements, on page 12](#)

Attach a Traffic Policy to an Interface

After the traffic class and the traffic policy are created, you must attach the traffic policy to interface, and specify the direction in which the policy should be applied.



Note Hierarchical policies are not supported.

When a policy-map is applied to an interface, the transmission rate counter of each class is not accurate. This is because the transmission rate counter is calculated based on the exponential decay filter.

Configuration Example

You have to accomplish the following to attach a traffic policy to an interface:

1. Creating a traffic class and the associated rules that match packets to the class (see [Create a Traffic Class, on page 14](#))
2. Creating a traffic policy that can be attached to one or more interfaces to specify a service policy (see [Create a Traffic Policy, on page 16](#))
3. Associating the traffic class with the traffic policy
4. Attaching the traffic policy to an interface, in the ingress or egress direction

```
Router# configure
Router(config)# interface fourHundredGigE 0/0/0/2
Router(config-int)# service-policy output strict-priority
Router(config-int)# commit
```

Running Configuration

```
/* Class-map configuration */

class-map match-any traffic-class-7
match traffic-class 7
end-class-map

!class-map match-any traffic-class-6
match traffic-class 6
end-class-map
```

```
class-map match-any traffic-class-5
match traffic-class 5
end-class-map

class-map match-any traffic-class-4
match traffic-class 4
end-class-map

class-map match-any traffic-class-3
match traffic-class 3

class-map match-any traffic-class-2
match traffic-class 2
end-class-map

class-map match-any traffic-class-1
match traffic-class 1
end-class-map

/* Traffic policy configuration */

policy-map test-shape-1
class traffic-class-1
  shape average percent 40
!

policy-map strict-priority
class tc7
  priority level 1
  queue-limit 75 mbytes
!
class tc6
  priority level 2
  queue-limit 75 mbytes
!
class tc5
  priority level 3
  queue-limit 75 mbytes
!
class tc4
  priority level 4
  queue-limit 75 mbytes
!
class tc3
  priority level 5
  queue-limit 75 mbytes
!
class tc2
  priority level 6
  queue-limit 75 mbytes
!
class tc1
  priority level 7
  queue-limit 75 mbytes
!
class class-default
  queue-limit 75 mbytes
!
end-policy-map

- - -
- - -
```

```

/* Attaching traffic policy to an interface in egress direction */
interface fourHundredGigE 0/0/0/2
 service-policy output strict-priority
 !

```

Verification

```
Router# #show qos int fourHundredGigE 0/0/0/2 output
```

NOTE:- Configured values are displayed within parentheses Interface FourHundredGigE0/0/0/2
ifh 0xf0001c0 -- output policy

```

NPU Id:                                0
Total number of classes:                8
Interface Bandwidth:                    400000000 kbps
Policy Name:                            strict-priority
VOQ Base:                               2400
Accounting Type:                        Layer1 (Include Layer 1 encapsulation and above)
-----
Level1 Class (HP1)                      = tc7
Egressq Queue ID                        = 2407 (HP1 queue)
Queue Max. BW.                          = no max (default)
TailDrop Threshold                       = 74999808 bytes / 2 ms (75 megabytes)
WRED not configured for this class

Level1 Class (HP2)                      = tc6
Egressq Queue ID                        = 2406 (HP2 queue)
Queue Max. BW.                          = no max (default)
TailDrop Threshold                       = 74999808 bytes / 2 ms (75 megabytes)
WRED not configured for this class

Level1 Class (HP3)                      = tc5
Egressq Queue ID                        = 2405 (HP3 queue)
Queue Max. BW.                          = no max (default)
TailDrop Threshold                       = 74999808 bytes / 2 ms (75 megabytes)
WRED not configured for this class

Level1 Class (HP4)                      = tc4
Egressq Queue ID                        = 2404 (HP4 queue)
Queue Max. BW.                          = no max (default)
TailDrop Threshold                       = 74999808 bytes / 2 ms (75 megabytes)
WRED not configured for this class

Level1 Class (HP5)                      = tc3
Egressq Queue ID                        = 2403 (HP5 queue)
Queue Max. BW.                          = no max (default)
TailDrop Threshold                       = 74999808 bytes / 2 ms (75 megabytes)
WRED not configured for this class

Level1 Class (HP6)                      = tc2
Egressq Queue ID                        = 2402 (HP6 queue)
Queue Max. BW.                          = no max (default)
TailDrop Threshold                       = 74999808 bytes / 2 ms (75 megabytes)
WRED not configured for this class

Level1 Class (HP7)                      = tc1
Egressq Queue ID                        = 2401 (HP7 queue)
Queue Max. BW.                          = no max (default)
TailDrop Threshold                       = 74999808 bytes / 2 ms (75 megabytes)
WRED not configured for this class

```

```

Level1 Class                = class-default
Egressq Queue ID           = 2400 (Default LP queue)
Queue Max. BW.             = no max (default)
Inverse Weight / Weight    = 1 / (BWR not configured)
TailDrop Threshold         = 74999808 bytes / 150 ms (75 megabytes)
WRED not configured for this class

```

```
!
```

Related Topics

- [Traffic Policy Elements, on page 15](#)
- [Traffic Class Elements, on page 12](#)

QoS Policy Propagation via BGP (QPPB)

Table 4: Feature History Table

Feature Name	Release Information	Feature Description
QoS Policy Propagation via BGP	Release 7.5.2	<p>You now have the ability to install a BGP route in the routing table with a QoS Group so that IP packets that match the route receive the QoS policies associated with the QoS group.</p> <p>This functionality enables convenient classification and marking when BGP is deployed, overcoming the administrative challenges of classifying based on ACLs.</p>

QoS Policy Propagation via Border Gateway Protocol (QPPB via BGP) is a mechanism that allows propagation of quality of service (QoS) policy and classification by the sending party that is based on the following:

- Access lists
- Community lists
- Autonomous system paths in the BGP

QPPB thus helps in classification that is based on the destination address instead of the source address.

QoS policies that differentiate between different types of traffic are defined for a single enterprise network. For instance, one enterprise may want to treat important web traffic, not-important web traffic, and all other data traffic as three different classes. And thereafter, use the different classes for the voice and video traffic.

Hence, QPPB overcomes:

- the administrative challenges of classifying that is based on ACLs.
- the administrative problems of just listing the networks that need premium services.

QPPB allows marking of packets that are based on QoS group value that is associated with a Border Gateway Protocol (BGP) route.

Benefits of QPPB

- QPPB provides an IP prefix-based QoS capability.
- Traffic to IP addresses that have specific IP prefixes can be prioritized above other IP addresses.
- IP prefixes of interest are tagged through the control plane that uses common BGP route-map techniques, including the community attribute.
- Traffic to the tagged BGP prefixes is then classified and prioritized via the data forwarding plane by using the IOS-XR MQC (Modular QoS CLI) mechanisms, such as remarking.
- QPPB provides the glue between the BGP control plane and the IP data forwarding plane in support of IP prefix-based QoS.
- BGP configuration within QPPB uses a table map to match specific prefixes learned through BGP neighbors, and then sets the router's local QoS Group variable that is maintained within the Forwarding Information Base (FIB) for those specific prefixes.

Guidelines and Limitations

- While configuring QPPB in the route policy for BGP, the range of the QoS group value that you can set using the **set qos-group** command is between 0 through 7.
- IP precedence-based QPPB isn't supported; only QoS group-based QPPB is supported.
- The order of precedence among ACL features (in ascending order):

With QoS-based group policy:

- Security ACL (0), BGP-FS (1), QPPB (2), Peering QoS (3)

With ACL-based policy:

- Security ACL (0), ACL based QoS (1), BGP-FS (2), QPPB (3)

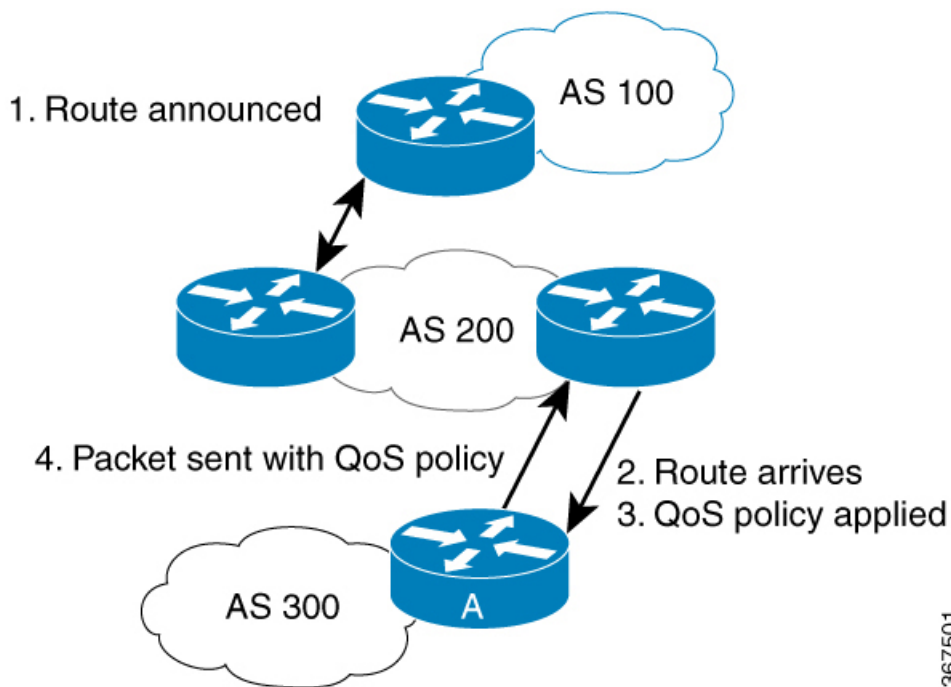
- You can configure QPPB on the following interface types:

- Physical
- Bundle
- Physical subinterface
- Bundle subinterface

- QPPB is not supported on Q100-based routers or line cards.
- QPPB overwrites the QoS Group value set by Security ACL.
- If you configure ACL-based Switched Port Analyzer (SPAN) on an interface, you can't enable QPPB on that interface. To enable QPPB, remove the ACL-based SPAN from the interface.
- If you configure QPPB on an interface, you can't enable ACL-based SPAN on that interface. To enable ACL-based SPAN, remove QPPB from the interface.

- Peering QoS overwrites the QoS Group value set by QPPB.
- Peering QoS and QPPB overwrite BGP FlowSpec actions except setting Drop and Policer. Currently, BGP FlowSpec sets DSCP and Meter.
- Remarking support is the same that QoS currently supports.
- Ingress policing support is the same that QoS currently supports.
- Egress policing is not supported.
- The router supports a subset of full QPPB options—only IP destination prefix mode on input policy is supported.

Figure 3: Sample Scenario



Router A learns routes from AS 200 and AS 100. QoS policy is applied to any ingress interface of Router A to match the defined route maps with destination prefixes of incoming packets. Matching packets on Router A to AS 200 or AS 100 are sent with the appropriate QoS policy from Router A.

BGP maintains a scalable database of destination prefixes, QPPB, by using BGP table maps. BGP adds the ability to map a qos-group value to desired IP destinations. These qos-group values are used in QoS policies applied locally on ingress interfaces. Whenever a packet that is bound for such destinations is encountered, the qos-group value matching that destination route looks up with work inside the policy classmap, and marks that packet for any configured policy actions.

Configuration Workflow

Use the following configuration workflow for QPPB:

- Define routing policy.

- Put routing policy at table-policy attach point under BGP.
- Define classmaps and ingress policy to use the qos-groups that are used in table-policy.
- Enable IPv4 or IPv6 QPPB configuration under the desired interfaces.

Define Routing Policy

A routing policy instructs the router to inspect routes, filter them, and potentially modify their attributes as they are accepted from a peer, advertised to a peer, or redistributed from one routing protocol to another.

The routing policy language (RPL) provides a language to express routing policy. You must set up destination prefixes either to match inline values or one of a set of values in a prefix set.

Example:

```
prefix-set prefix-list-v4
 70.1.1.1,
 70.2.1.0/24,
 70.2.2.0/24 ge 28,
 70.2.3.0/24 le 28
end-set
prefix-set prefix-list-v6
 2001:300::2,
 2003:200::3
end-set

route-policy qppb1
  if destination in (60.60.0.2/24) then
    set qos-group 5
  elseif destination in prefix-list-v4 then
    set qos-group 4
  else
    set qos-group 1
  pass
endif
end-policy
route-policy qppb2
  if destination in prefix-list-v6 then
    set qos-group 5
  elseif destination in (2001:300::2) then
    set qos-group 4
  else
    set qos-group 1
  pass
endif
end-policy
```

Put Routing Policy at table-policy Attach Point Under BGP

The table-policy attach point permits the routing policy to perform actions on each route as they are installed into the RIB routing table. QPPB uses this attachment point to intercept all routes as they are received from peers. Ultimately the RIB will update the FIB in the hardware forwarding plane to store destination prefix routing entries, and in cases where table policy matches a destination prefix, the qos-group value is also stored with the destination prefix entry for use in the forwarding plane.

Example:

```
router bgp 900
 [vrf <name>]
  bgp router-id 22.22.22.22
  address-family ipv4 unicast
```

```

    table-policy qppb1
address-family ipv6 unicast
    table-policy qppb2
neighbor 30.2.2.1
    remote-as 500
    address-family ipv4 unicast
        route-policy pass in
        route-policy pass out
    address-family ipv6 unicast
        route-policy pass in
        route-policy pass out

```

Ingress QoS and IPv4 or IPv6 BGP Configuration

QPPB is enabled per interface and individually for IPv4 and IPv6. An ingress policy matches on the QoS groups marked by QPPB and takes the required action.

If a packet is destined for a destination prefix on which BGP route policy has stored a **qos-group**, but it ingresses on an interface on which QPPB is not enabled, it would not be remarked with **qos-group**.

Example:

```

class-map match-any qos-group5
    match qos-group 5
end-class-map

class-map match-any qos-group4
    match qos-group 4
end-class-map

policy-map ingress-marker-pol
    class qos-group5
        set precedence 0
        set discard-class 0
        set traffic-class 1

    class qos-group4
        set precedence 1
        set discard-class 1
        set traffic-class 2
    class class-default

end-policy-map

Interface hun 0/0/0/0
    [vrf vrfA]
    ipv4 address 25.1.1.1/24
    ipv6 address 2001:db8:a0b:12f0::1/64
    ipv4 bgp policy propagation input qos-group destination
    ipv6 bgp policy propagation input qos-group destination
    service-policy input ingress-marker-pol

```

Egress Interface Configuration

The traffic-class set on ingress has no existence outside the device. Also, traffic-class is not a part of any packet header but is associated internal context data on relevant packets. It can be used as a match criteria in an egress policy to set up various fields on the outgoing packet or shape flows.

```

class-map match-any level1
    match traffic-class 1
end-class-map

```

```
class-map match-any level2
  match traffic-class 2
end-class-map

policy-map output-pol
  class level1
    bandwidth percent 50
  class level2
    bandwidth percent 20
    queue-limit 50 ms
end-policy-map

interface hun 0/5/0/0
  ipv4 address 30.1.1.1/24
  ipv6 address 2001:da8:b0a:12f0::1/64
  service-policy output output-pol
```




CHAPTER 4

Mark Packets to Change Priority Settings

- [Packet Marking Overview](#), on page 27
- [Class-based Unconditional Packet Marking Feature and Benefits](#), on page 29
- [Configure Class-based Unconditional Packet Marking](#), on page 30
- [Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic](#), on page 32
- [Class-based Unconditional Packet Marking: Examples](#), on page 35
- [IP Precedence Compared to IP DSCP Marking](#), on page 40
- [In-Place Policy Modification](#), on page 41

Packet Marking Overview

You can use packet marking in input policy maps to set or modify the attributes for traffic belonging to a specific class. For example, you can change the CoS value in a class or set IP DSCP or IP precedence values for a specific type of traffic. These new values are then used to determine how the traffic should be treated.



Note From Cisco IOS XR Release 7.2.12 onwards, support for marking packets on Layer 2 transport interfaces is the same as the support for marking on Layer 3 interfaces. However, this support applies only to the main interface (physical and bundle interfaces), and not on the sub-interfaces.

Default Marking

When an ingress or egress interface adds VLAN tags or MPLS labels, it requires a default value for the class of service and EXP values that go into those tags and labels.

On the router, one ingress default QoS mapping profile and one egress default QoS mapping profile are created and configured per device during initialization.

QoS Behavior for Generic Routing Encapsulation (GRE) Tunnels

Table 5: Feature History Table

Feature Name	Release Information	Feature Description
QoS Behavior for Generic Routing Encapsulation (GRE) Tunnels: Default Marking	Release 7.3.1	With the support for GRE encapsulation and decapsulation tunnel interfaces, there are some important updates to QoS behavior for GRE tunnels. These updates are applicable for default packet marking and involve Type of Service (ToS) and MPLS experimental bits.

GRE Encapsulation

If you do not configure Type of Service (ToS), the outer IP precedence value or the differentiated services code point (DSCP) value is copied from the inner IP header. If you configure ToS, the outer IP precedence value or DSCP value is as per the ToS configuration.

GRE Decapsulation

During decapsulation, the MPLS experimental bits (EXP) are derived from the outer IP packet.

For more information on GRE tunnels, see the *Interfaces Configuration Guide for Cisco 8000 Series Routers, IOS XR Release 7.3.x*.

Packet Marking

The packet marking feature, also called explicit marking, provides users with a means to differentiate packets based on the designated markings. The router supports ingress and egress packet marking.

Supported Packet Marking Operations

This table shows the supported packet marking operations.

Table 6: Supported Packet Marking Operations for Egress and Ingress

Supported Mark Types	Range	Support for Unconditional Marking	Support for Conditional Marking
set discard-class	0-1	ingress	No
set dscp	0-63	ingress, egress	No
set mpls experimental topmost	0-7	ingress, egress	No
set precedence	0-7	ingress, egress	No
set qos-group	0-7	ingress	No

QoS Behavior for Generic Routing Encapsulation (GRE) Tunnels

Table 7: Feature History Table

Feature Name	Release Information	Feature Description
QoS Behavior for Generic Routing Encapsulation (GRE) Tunnels: Explicit Marking	Release 7.3.1	With the support for GRE encapsulation and decapsulation tunnel interfaces, there are some important updates to QoS behavior for GRE tunnels. These updates are applicable for explicit packet marking and involve QoS behavior during ingress and egress.

GRE Encapsulation

During encapsulation of IPv4/IPv6 payload inside the GRE header, QoS behavior is as follows:

- Ingress: QoS supports classification on the payload Layer 3 fields or EXP and remarking payload IP header DSCP.
- Egress: QoS supports setting outer GRE IP header DSCP. It doesn't overwrite the Tunnel Type of Service (ToS) configuration and doesn't remark GRE IP header DCSP.

GRE Decapsulation

During decapsulation of the outer GRE header (during which the inner IPv4/IPv6/MPLS payload is forwarded to the next-hop router), QoS behavior is as follows:

- Ingress: QoS supports classification on Layer 3 fields of outer GRE using the **set qos-group** command. Setting DSCP on the ingress interface sets DSCP for the inner headers.
- Egress: QoS supports classification using **qos-group** to set DSCP or EXP for egress packets.

For more information on GRE tunnels, see the *Interfaces Configuration Guide for Cisco 8000 Series Routers, IOS XR Release 7.3.x*.

Class-based Unconditional Packet Marking Feature and Benefits

The packet marking feature allows you to partition your network into multiple priority levels or classes of service, as follows:

- Use QoS unconditional packet marking to set the IP precedence or IP DSCP values for packets entering the network. Routers within your network can then use the newly marked IP precedence values to determine how the traffic should be treated.

On ingress direction, after matching the traffic based on either the IP Precedence or DSCP value, you can set it to a particular discard-class. Weighted random early detection (WRED), a congestion avoidance technique, thereby uses discard-class values to determine the probability that a packet is dropped.

- Use QoS unconditional packet marking to assign MPLS packets to a QoS group. The router uses the QoS group to determine how to prioritize packets for transmission. To set the QoS group identifier on MPLS packets, use the **set qos-group** command in policy map class configuration mode.



Note Setting the QoS group identifier does not automatically prioritize the packets for transmission. You must first configure an egress policy that uses the QoS group.

- Mark Multiprotocol Label Switching (MPLS) packets by setting the EXP bits within the imposed or topmost label.
- Mark packets by setting the value of the *qos-group* argument.
- Mark packets by setting the value of the *discard-class* argument.



Note *qos-group* and *discard-class* are variables internal to the router, and are not transmitted.

The configuration task is described in the [Configure Class-based Unconditional Packet Marking, on page 30](#).

Configure Class-based Unconditional Packet Marking

This configuration task explains how to configure the following class-based, unconditional packet marking features on your router:

- IP precedence value
- IP DSCP value
- QoS group value (ingress only)
- CoS value (egress only on Layer 3 subinterfaces)
- MPLS experimental value
- Discard class



Note IPv4 and IPv6 QoS actions applied to MPLS tagged packets are not supported. The configuration is accepted, but no action is taken.

Configuration Example

Follow these steps to configure unconditional packet marking features on your router.

1. Create or modify a policy map that can be attached to one or more interfaces to specify a service policy and enter the policy map configuration mode.
2. Configure an interface and enter the interface configuration mode.

3. Attach a policy map to an input or output interface to be used as the service policy for that interface.

Configuration Example

```
router# configure
router(config)# interface hundredGigE 0/0/0/24
router(config-pmap)# policy-map policy1
Router(config-int)# commit
```

Running Configuration

```
router(config)# policy-map policy1

class-map match-any class1
  match protocol ipv4
end-class-map
!
!
policy-map policy1
  class class1
    set precedence 1
  !
  class class-default
  !
end-policy-map
!
interface HundredGigE0/0/0/24
  service-policy input policy1
!
!
```

Verification

Run this command to display policy configuration information for all classes configured for all service policies on the specified interface.

```
router# show run interface hundredGigE 0/0/0/24
```

Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic

Table 8: Feature History Table

Feature Name	Release Information	Feature Description
Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic	Release 7.5.4	<p>We've provided you with more granular control for IP to MPLS traffic, so you can collect billing information for your customers upstream at the MPLS core and correlate them with subsequent next-hop IP nodes. This helps you accurately bill your customers as they pay for specific services based on their network requirements.</p> <p>Such control is possible because we've enabled the remarking of IP DSCP values in the inner IP header—that remain even after the POP tag removes the outer header—along with the existing ability to remark MPLS EXP values on the outer MPLS header.</p> <p>There's no action required for you to enable this functionality, and there are no changes to command options.</p>

In earlier releases, when you remarked the DSCP value in the ingress policy for IP to MPLS traffic, here's what happened post the MPLS imposition:

- The inner IP header didn't receive the remarked DSCP value, meaning that the QoS action you required for the IP packet across the network didn't take place.
- The outer header received the remarked value for EXP bits based on your remarking action.

From Cisco IOS XR Release 7.5.4, for IP traffic flowing across an MPLS core, here's the QoS behavior when you remark or set the DSCP value in the ingress policy:

- At the ingress, when the MPLS header is imposed on an IP packet, the inner IP header receives the DSCP value you remarked in the ingress policy. In other words, the IP DSCP value in the inner header is remarked.
- As in earlier releases, the top label in the MPLS header receives the remarked value for EXP bits based on your remarking action. This value for the EXP bits is derived from the three most significant bits of the remarked DSCP value.

Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic: Benefits

With the ability to remark IP DSCP and MPLS EXP values at the ingress, you have granular control of QoS actions across hops in your network. You can now collect customer-specific billing information upstream at the MPLS core and collate it with the traffic service provided at various IP locations.

Mark IP DSCP and MPLS EXP Values for IP to MPLS Traffic: Guidelines and Limitations

- This functionality is available only on Cisco Silicon One Q200-based fixed and distributed systems.
- You can use a Cisco Silicon One Q100-based line card for distributed systems, but the egress line card must be Cisco Silicon One Q200-based. Else, the feature won't work.

Configure IP DSCP and MPLS EXP Values for IP to MPLS Traffic

Configuration Example

In this example, you set the DSCP value (or assured forwarding per-hop behavior [PHB]) as AF41 for the policy map named **INGRESS_REMARK** at the ingress. At the ingress PE router, the configuration would be:

```
Router(config)#class-map AF31
Router(config-cmap)#match dscp AF31
Router(config-cmap)#end-class-map
Router(config)#
Router(config)#policy-map INGRESS_REMARK
Router(config-pmap)#class AF31
Router(config-pmap-c)#set dscp AF41
Router(config-pmap-c)#end-policy-map
Router(config)#
Router(config)#interface FourHundredGigE0/0/0/0
Router(config-if)#service-policy input INGRESS_REMARK
Router(config-if)#commit
```

At the egress PE router, your configuration would be:

```
Router(config)#class-map AF31
Router(config-cmap)#match dscp AF31
Router(config-cmap)#end-class-map
Router(config)#
Router(config)#class-map AF41
Router(config-cmap)#match dscp AF41
Router#end-class-map
Router(config)#
Router(config)#policy-map INGRESS_CLASSIFY
Router(config-pmap)#class AF31
Router(config-pmap-c)#set traffic-class 3
Router(config-pmap-c)#exit
Router(config-pmap)#
Router(config-pmap)#class AF41
Router(config-pmap-c)#set traffic-class 4
Router(config-pmap-c)#exit
Router(config-pmap)#end-policy-map
```

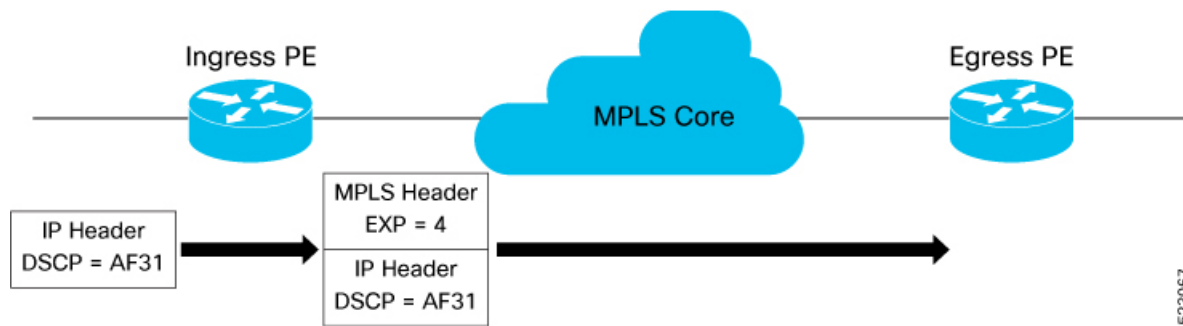
```

Router(config)#
Router(config)#int FourHundredGigE0/0/0/1
Router(config-if)#service-policy input INGRESS_CLASSIFY
Router(config-if)#commit

```

Verification

Figure 4: MPLS QoS Behavior Before Cisco IOS XR Release 7.5.4



Before this release, even though you configured the new DSCP value as AF41 for the policy map named **INGRESS_REMARK** at the ingress, the output for the **show policy-map interface** showed that the assured forwarding marking wasn't applied for the PHB. Thus, the statistics counter displayed zero for Class AF41 at the egress PE router.

```
Router#show policy-map interface FourHundredGigE0/0/0/1 input
```

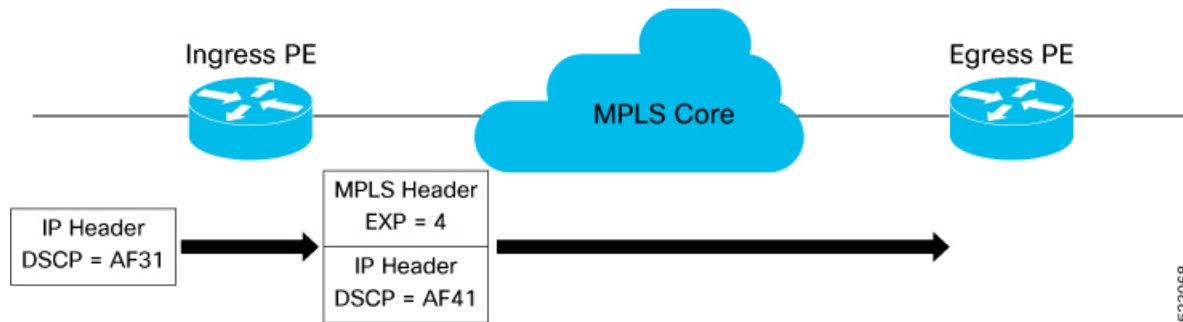
```
FourHundredGigE0/0/0/1 input: INGRESS_CLASSIFY
```

```

Class AF31
  Classification statistics      (packets/bytes)      (rate - kbps)
  Matched                      :      9813350/13738690000      936847
  Transmitted                   :      9813350/13738690000      936847
  Total Dropped                 :                0/0                0
Class AF41
  Classification statistics      (packets/bytes)      (rate - kbps)
  Matched                      :                0/0                0
  Transmitted                   :                0/0                0
  Total Dropped                 :                0/0                0
Class class-default
  Classification statistics      (packets/bytes)      (rate - kbps)
  Matched                      :                0/0                0
  Transmitted                   :                0/0                0
  Total Dropped                 :                0/0                0

```

Figure 5: MPLS QoS Behavior From Cisco IOS XR Release 7.5.4



With the ability to remark IP DSCP values in the inner IP header, the **show policy-map interface** command now displays traffic statistics for Class AF41 at the egress PE router. These statistics confirm that the requisite QoS action—which is, DSCP remark—has taken place at the ingress PE router.

```
Router#show policy-map interface FourHundredGigE0/0/0/1 input

FourHundredGigE0/0/0/1 input: INGRESS_CLASSIFY

Class AF31
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          :                   0/0                0
  Transmitted                       :                   0/0                0
  Total Dropped                     :                   0/0                0
Class AF41
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          :          9813350/13738690000    936847
  Transmitted                       :          9813350/13738690000    936847
  Total Dropped                     :                   0/0                0
Class class-default
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          :                   0/0                0
  Transmitted                       :                   0/0                0
  Total Dropped                     :                   0/0                0
```

Class-based Unconditional Packet Marking: Examples

These are typical examples for class-based unconditional packet marking.

IP Precedence Marking Configuration: Example

In this example, a service policy called *policy1* is created. This service policy is associated to a previously defined class map called *class1* through the use of the **class** command, and then the service policy is attached to the output HundredGigE interface 0/7/0/1. The IP precedence bit in the ToS byte is set to 1:

```
policy-map policy1
  class class1
    set precedence 1
!
interface HundredGigE 0/7/0/1
  service-policy output policy1
```

IP DSCP Marking Configuration: Example

In this example, a service policy called *policy1* is created. This service policy is associated to a previously defined class map through the use of the **class** command. In this example, it is assumed that a class map called *class1* was previously configured and new class map called *class2* is created.

In this example, the IP DSCP value in the ToS byte is set to 5:

```
policy-map policy1
  class class1
    set dscp 5
```

```
class class2
  set dscp ef
```

After you configure the settings shown for voice packets at the edge, all intermediate routers are configured to provide low-latency treatment to the voice packets, as follows:

```
class-map voice
  match dscp ef
policy-map qos-policy
  class voice
    priority level 1
    police rate percent 10
```

QoS Group Marking Configuration: Example

In this example, a service policy called *policy1* is created. This service policy is associated to a class map called *class1* through the use of the **class** command, and then the service policy is attached in the input direction on a HundredGigE 0/7/0/1. The qos-group value is set to 1.

```
class-map match-any class1
  match protocol ipv4
  match access-group ipv4 101

policy-map policy1
  class class1
    set qos-group 1
  !
interface HundredGigE 0/7/0/1
  service-policy input policy1
```



Note The **set qos-group** command is supported only on an ingress policy.

CoS Marking Configuration: Example

In this example, a service policy called *policy1* is created. This service policy is associated to a class map called *class1* through the use of the **class** command, and then the service policy is attached in the output direction on a HundredGigE 0/7/0/1.100. The IEEE 802.1p (CoS) bits in the Layer 2 header are set to 1.

```
class-map match-any class1
  match protocol ipv4
  match access-group ipv4 101

policy-map policy1
  class class1
    set cos 1
  !

interface HundredGigE 0/7/0/1.100
  service-policy input policy1
```

MPLS Experimental Bit Imposition Marking Configuration: Example

In this example, a service policy called *policy1* is created. This service policy is associated to a class map called *class1* through the use of the **class** command, and then the service policy is attached in the input direction on a HundredGigE 0/7/0/1. The MPLS EXP bits of all imposed labels are set to 1.

```
class-map match-any class1
  match protocol ipv4
  match access-group ipv4 101

policy-map policy1
  class class1
    set mpls exp imposition 1
  !
interface HundredGigE 0/7/0/1
  service-policy input policy1
```



Note The **set mpls exp imposition** command is supported only on an ingress policy.

MPLS Experimental Topmost Marking Configuration: Example

In this example, a service policy called *policy1* is created. This service policy is associated to a class map called *class1* through the use of the **class** command, and then the service policy is attached in the output direction on a HundredGigE 0/7/0/1. The MPLS EXP bits on the TOPMOST label are set to 1:

```
class-map match-any class1
  match mpls exp topmost 2

policy-map policy1
  class class1
    set mpls exp topmost 1
  !
interface HundredGigE 0/7/0/1
  service-policy output policy1
```

Queuing and Marking Policies: Examples

These examples explain how queuing and marking can be applied to an interface.

Ingress Traffic Classification and QoS Mapping Configuration for Input Interfaces

In this example, a policy map named "INGRESS_L3" is configured, specifically for classifying and marking incoming traffic.

```
RP/0/RP0/CPU0:Router#show policy-map pmap-name INGRESS_L3 detail

class-map match-all CONTROL_PLANE
match dscp cs7
end-class-map
!
class-map match-all VOIP
match dscp cs6
end-class-map
```

```

!
class-map match-all VIDEO_STREAM
match dscp cs5
end-class-map
!
class-map match-all TRANSACTIONAL_DATA
match dscp cs4
end-class-map
!
class-map match-all DB_SYNC
match dscp cs3
end-class-map
!
class-map match-all BULK_DATA
match dscp cs2
end-class-map
!
class-map match-all SCAVENGER
match dscp cs1
end-class-map
!
policy-map INGRESS_L3
class CONTROL_PLANE
set traffic-class 7
set qos-group 27
!
class VOIP
set traffic-class 6
set qos-group 26
!
class VIDEO_STREAM
set traffic-class 5
set qos-group 25
!
class TRANSACTIONAL_DATA
set traffic-class 4
set qos-group 24
!
class DB_SYNC
set traffic-class 3
set qos-group 23
!
class BULK_DATA
set traffic-class 2
set qos-group 22
!
class SCAVENGER
set traffic-class 1
set qos-group 21
!
class class-default
!
end-policy-map
!

```

Egress Queueing Policy Configuration

In this example, a policy map named "EGRESS_QOS_QUEUEING" is configured for the egress (outgoing) direction.

```
RP/0/RP0/CPU0:Router#show policy-map pmap-name EGRESS_QOS_QUEUEING detail
```

```
class-map match-any TC7_CONTROL_PLANE
match traffic-class 7
```



```

    end-class-map
    !
    class-map match-any TC6_VOIP
    match traffic-class 6
    end-class-map
    !
    class-map match-any TC5_VIDEO_STREAM
    match traffic-class 5
    end-class-map
    !
    class-map match-any TC4_TRANSACTIONAL_DATA
    match traffic-class 4
    end-class-map
    !
    class-map match-any TC3_DB_SYNC
    match traffic-class 3
    end-class-map
    !
    class-map match-any TC2_BULK_DATA
    match traffic-class 2
    end-class-map
    !
    class-map match-any TC1_SCAVENGER
    match traffic-class 1
    end-class-map
    !
    policy-map EGRESS_QOS_QUEUEING
    class TC7_CONTROL_PLANE
    priority level 1
    !
    class TC6_VOIP
    priority level 2
    !
    class TC5_VIDEO_STREAM
    bandwidth remaining ratio 20
    random-detect 50 mbytes 100 mbytes
    !
    class TC4_TRANSACTIONAL_DATA
    bandwidth remaining ratio 20
    queue-limit 75 mbytes
    !
    class TC3_DB_SYNC
    bandwidth remaining ratio 10
    !
    class TC2_BULK_DATA
    bandwidth remaining ratio 10
    !
    class TC1_SCAVENGER
    bandwidth remaining ratio 5
    !
    class class-default
    !
    end-policy-map
    !

```

Egress Marking Policy Configuration

In this example, a policy map named "EGRESS_REMARK" is configured for the egress (outgoing) direction.

```
RP/0/RP0/CPU0:Router#show policy-map pmap-name EGRESS_REMARK detail
```

```

class-map match-any QG24
match qos-group 24
end-class-map

```

```

!
class-map match-any QG23
match qos-group 23
end-class-map
!
class-map match-any QG22
match qos-group 22
end-class-map
!
policy-map EGRESS_REMARK
class QG24
set dscp af42
!
class QG23
set dscp af33
!
class QG22
set dscp af21
!
class class-default
!
end-policy-map
!

```

Output Interface Configuration with Queuing Policy and Remark Policy Applied

```

interface FourHundredGigE0/2/0/29
mtu 10000
service-policy output EGRESS_REMARK
service-policy output EGRESS_QOS_QUEUEING
ipv4 address 18.29.0.1 255.255.255.0
ipv6 address 18:1d::1/96
!

```

IP Precedence Compared to IP DSCP Marking

If you need to mark packets in your network and all your devices support IP DSCP marking, use the IP DSCP marking to mark your packets because the IP DSCP markings provide more unconditional packet marking options. If marking by IP DSCP is undesirable, however, or if you are unsure if the devices in your network support IP DSCP values, use the IP precedence value to mark your packets. The IP precedence value is likely to be supported by all devices in the network.

You can set up to 8 different IP precedence markings and 64 different IP DSCP markings.

Configure DSCP CS7 (Precedence 7)

See the following example to configure options in DSCP for a particular source address in IPv4 packets.

Configuration Example

```

policy-map policy1
class class1
set dscp cs7
!

```

In-Place Policy Modification

The In-Place policy modification feature allows you to modify a QoS policy even when the QoS policy is attached to one or more interfaces. A modified policy is subjected to the same checks that a new policy is subject to when it is bound to an interface. If the policy-modification is successful, the modified policy takes effect on all the interfaces to which the policy is attached. However, if the policy modification fails on any one of the interfaces, an automatic rollback is initiated to ensure that the pre-modification policy is in effect on all the interfaces.

You can also modify any class map used in the policy map. The changes made to the class map take effect on all the interfaces to which the policy is attached.



Note

- The QoS statistics for the policy that is attached to an interface are lost (reset to 0) when the policy is modified.
 - When a QoS policy attached to an interface is modified, there might not be any policy in effect on the interfaces in which the modified policy is used for a short period of time.
 - An in-place modification of an ACL does not reset the policy-map statistics counter.
-

Verification

If unrecoverable errors occur during in-place policy modification, the policy is put into an inconsistent state on target interfaces. No new configuration is possible until the configuration session is unblocked. It is recommended to remove the policy from the interface, check the modified policy and then re-apply accordingly.

Recommendations for Using In-Place Policy Modification

For a short period of time while a QoS policy is being modified, there might not be any policy in effect on the interfaces in which the modified policy is used. For this reason, modify QoS policies that affect the fewest number of interfaces at a time. Use the **show policy-map targets** command to identify the number of interfaces that will be affected during policy map modification.



CHAPTER 5

Congestion Avoidance

- [Congestion Avoidance](#), on page 43
- [Queuing Modes](#), on page 43
- [Congestion Avoidance in VOQ](#), on page 44
- [Equitable Traffic Flow Using Fair VOQ](#), on page 48
- [Modular QoS Congestion Avoidance](#), on page 54
- [Tail Drop and the FIFO Queue](#), on page 54
- [Random Early Detection and TCP](#), on page 56
- [Explicit Congestion Notification](#), on page 58

Congestion Avoidance

Queuing provides a way to temporarily store data when the received rate of data is larger than what can be sent. Managing the queues and buffers is the primary goal of congestion avoidance. As a queue starts to fill up with data, it is important to try to make sure that the available memory in the ASIC/NPU does not fill up completely. If this happens, subsequent packets coming into the port are dropped, irrespective of the priority that they received. This could have a detrimental effect on the performance of critical applications. For this reason, congestion avoidance techniques are used to reduce the risk of a queue from filling up the memory completely and starving non-congested queues for memory. Queue thresholds are used to trigger a drop when certain levels of occupancy are exceeded.

Scheduling is the QoS mechanism that is used to empty the queues of data and send the data onward to its destination.

Shaping is the act of buffering traffic within a port or queue until it is able to be scheduled. Shaping smoothens traffic, making traffic flows much more predictable. It helps ensure that each transmit queue is limited to a maximum rate of traffic.

Queuing Modes

Two network queuing modes are supported for network interface queuing: the default mode of 8xVOQ (virtual output queuing) and 4xVOQ. To change the mode from one to another requires that you must first reload all line cards in the system.

In the 8xVOQ mode, eight VoQs and their associated resources are allocated for each interface. These queues are allocated regardless of the exact policy configuration on that interface. This mode supports a separate VOQ for each of the eight internal traffic classes.

In the 4xVOQ mode, four VoQs and their associated resources are allocated to each interface, and these queues are allocated regardless of the exact policy applied. In this mode the system supports twice the number of logical interfaces, but the eight traffic classes must be mapped down by configuration to four VoQs, not eight VoQs.



Note From Cisco IOS XR Release 7.2.12 onwards, all the queuing features that are supported on Layer 3 interfaces are also supported on Layer 2 interfaces. However, these features apply only to the main interface (physical and bundle interfaces), and not on the sub-interfaces.

Main Interface Queuing Policy

The main interface default queues are created as part of the main interface creation.

When you apply a queuing policy to the main interface, it will override the default queuing and scheduling parameters for the traffic classes you have configured.

In the 8xVOQ mode, a P1+P2+6PN hierarchy is used for the main interface queues (default queuing and scheduling). The default queues are used for all traffic to the main interface and traffic to any sub-interface without a queuing policy applied. The control/protocol traffic uses traffic class 7 (TC7), priority 1 (P1) to avoid drops during congestion.

Sub-Interface Queuing Policy

Each sub-interface supports up to three policies: an ingress policy, an egress marking policy, and an egress queuing policy. To create and configure a separate set of VoQs for a sub-interface, apply a queuing policy on that sub-interface. When you remove the sub-interface queuing policy, the associated VoQs are freed and the sub-interface traffic reverts to using the main interface VoQs.

Congestion Avoidance in VOQ

Congestion avoidance within a VOQ block is done by applying a congestion management profile to a VOQ. This profile defines the admission criteria and checks performed at the enqueue time. Under normal traffic conditions the packet is enqueued into the Shared Memory System (SMS) buffers. (The shared memory system is the primary packet storage area.) If the SMS VOQ is congested beyond a set threshold, the VOQ is moved to the external High Band Memory (HBM) block. When the HBM queue drains, it is returned to the on-chip SMS. The queue size in HBM is adaptive and decreases when the total HBM usage is high.



Note Random Early Detect (RED) is available only for VOQs in HBM. The hardware does not support Weighted Random Early Detect (WRED).

Sharing of VOQ Statistics Counters

Every network processor on the router has multiple slices (or pipelines), and every slice has a set of VOQs associated with every interface on the router. To maintain counters at high packet rates, two sets of counters are associated with each interface on each network slice. As an example, consider a device with six slices (12 interfaces), each with 24,000 VOQs, where you want both transmitted and dropped events counted. In this scenario, you would require $12 \times 24,000 \times 2 = 5,76,000$ counters, which alone exceeds the counter capacity of the device.

It is to mitigate such a scenario that the router supports configurable sharing of VOQ counters. You can configure the sharing such that a counter is shared by {1,2,4,8} VOQs.

Each set of VoQs sharing counters has two counters that measure:

- Enqueued packets count in packets and bytes units.
- Dropped packets count in packets and bytes units.

For the feature to take effect:

- Delete egress queuing policy-map configuration from all interfaces.
- Run the command **# reload location all** to reload all the nodes on your router.

Configuring Sharing of VOQ Statistics Counters

To configure VOQs sharing counters, use the `#hw-module profile stats voqs-sharing-counters` and specify the number of VOQ counters for each queue.

```
RP/0/RP0/CPU0:ios(config)#hw-module profile stats ?
  voqs-sharing-counters  Configure number of voqs (1, 2, 4) sharing counters
RP/0/RP0/CPU0:ios(config)#hw-module profile stats voqs-sharing-counters ?
  1  Counter for each queue
  2  2 Queues share counters
  4  4 Queues share counters
RP/0/RP0/CPU0:ios(config)#hw-module profile stats voqs-sharing-counters 1
RP/0/RP0/CPU0:ios(config)#hw-module profile stats voqs-sharing-counters 2
RP/0/RP0/CPU0:ios(config)#commit
RP/0/RP0/CPU0:ios#reload location all
```

Running Configuration

```
RP/0/RP0/CPU0:ios#show run | in hw-mod
Mon Feb 10 13:57:35.296 UTC
Building configuration...
hw-module profile stats voqs-sharing-counters 2
RP/0/RP0/CPU0:ios#
```

Verification

```
RP/0/RP0/CPU0:ios#show controllers npu stats voq ingress interface hundredGigE 0/0/0/16
instance all location 0/RP0/CPU0
Mon Feb 10 13:58:26.661 UTC
```

```
Interface Name      = Hu0/0/0/16
Interface Handle    = f0001b0
Location            = 0/RP0/CPU0
Asic Instance       = 0
VOQ Base            = 10288
```

```

Port Speed(kbps) = 100000000
Local Port = local
VOQ Mode = 8
Shared Counter Mode = 2
      ReceivedPkts   ReceivedBytes   DroppedPkts   DroppedBytes
-----
TC_{0,1} = 114023724   39908275541   113945980     39881093000
TC_{2,3} = 194969733   68239406550   196612981     68814543350
TC_{4,5} = 139949276   69388697075   139811376     67907466750
TC_{6,7} = 194988538   68242491778   196612926     68814524100

```

Related Commands hw-module profile stats voqs-sharing-counters

Dual Queue Limit

The dual queue limit option is added to **queue-limit** command on the CLI of your router and displays as **discard-class**. What the **discard-class** option does is give you the flexibility to configure two queue limits on a single policy map—one for high-priority traffic and the other for low-priority traffic. This option ensures that the high priority traffic flow continues unaffected (up to the derived threshold from the **discard-class 0** queue-limit) while the low-priority traffic continues up to the lower threshold (per **discard-class 1** queue-limit).

Tell Me More

You can configure the two queue limits per these details:

- One for flow that you mark as **discard-class 0** (higher priority) on ingress via ingress-policy.
- second, for flow that you mark as **discard-class 1** (lower priority) on ingress via ingress policy.

The **discard-class 1** flow (for low-priority traffic) begins to drop when the queue length hits the size limit that you configured for discard-class 1. Conversely, the flow for **discard-class 1** stops dropping when queue-length falls below its configured value.

As an example, consider this configuration:

```

policy-map egress_pol_dql
class tc7
  queue-limit discard-class 0 100 mbytes
  queue-limit discard-class 1 50 mbytes
  priority level 1
!
class class-default
  bandwidth remaining ratio 1
!
end-policy-map
!

```

Also consider the verification:

```

RP/0/RP0/CPU0:ios#
RP/0/RP0/CPU0:ios#show qos interface hundredGigE 0/0/0/30 output
NOTE:- Configured values are displayed within parentheses
Interface HundredGigE0/0/0/30 ifh 0xf000210 -- output policy
NPU Id: 0
Total number of classes: 2
Interface Bandwidth: 100000000 kbps
Policy Name: egress_pol_dql
VOQ Base: 464
Accounting Type: Layer1 (Include Layer 1 encapsulation and above)
VOQ Mode: 8

```



```

Shared Counter Mode:          1
-----
Levell Class (HP1)           =   tc7
Egressq Queue ID             =   471 (HP1 queue)
Queue Max. BW.               =   no max (default)
Discard Class 1 Threshold     =   25165824 bytes / 2 ms (50 mbytes)
Discard Class 0 Threshold     =   75497472 bytes / 5 ms (100 mbytes)
WRED not configured for this class

Levell Class                 =   class-default
Egressq Queue ID             =   464 (Default LP queue)
Queue Max. BW.               =   no max (default)
Inverse Weight / Weight      =   1 / (1)
TailDrop Threshold           =   749568 bytes / 6 ms (default)
WRED not configured for this class

```

In the preceding example, there are two traffic flows that are marked as **discard-class 0** (higher priority) and **discard-class 1** (lower priority).

As long as the queue length of the two flows remains below 25165824 bytes (the threshold for **discard-class 1**), packets from both flows continue without any drops. When the queue length reaches 25165824 bytes, **discard-class 1** packets are not enqueued, ensuring all remaining bandwidth is used for the higher priority flow (**discard-class 0**).

The higher priority flow drops only when the queue length reaches 75497472 bytes.



Note

- This option protects the high-priority traffic from loss due to congestion, but not necessarily from latency due to congestion.
- These thresholds are derived from hardware-specific queue regions.

Restrictions

Ensure that you read these restrictions about the dual queue limit option.

- Both the queue-limits must use the same unit of measurement.
- The queue limit for **discard-class 0** must always be greater than that for **discard-class 1**.
- When the discard-class option is not used to configure the queue-limit, packets marked with **discard-class 0** and **discard-class 1** have the same queue-limit; in other words, they receive identical treatment.
- A queue-limit that is configured with only **discard-class 0** or **discard-class 1** is rejected.

Equitable Traffic Flow Using Fair VOQ

Table 9: Feature History Table

Feature Name	Release Information	Feature Description
Equitable Traffic Flow Using Fair VOQ	Release 7.3.3	<p>Configuring this feature ensures that ingress traffic from various source ports on every network slice of an NPU is assigned a unique virtual output queue (VOQ) for every source port and destination port pair. This action ensures that the bandwidth available at the destination port for a given traffic class is distributed equally to all source ports requesting bandwidth.</p> <p>In earlier releases, the traffic wasn't distributed equitably because each slice wasn't given its fair share of the output queue bandwidth.</p> <p>This feature introduces the fair-4 and fair-8 keywords in the hw-module profile qos voq-mode command.</p>
Equitable Traffic Flow Using Fair VOQ on Cisco 8201-32FH Routers	Release 7.5.2	<p>You can now ensure that the bandwidth available at the destination port for a given traffic class is distributed equally to all source ports requesting bandwidth on Cisco 8201-32FH routers.</p>

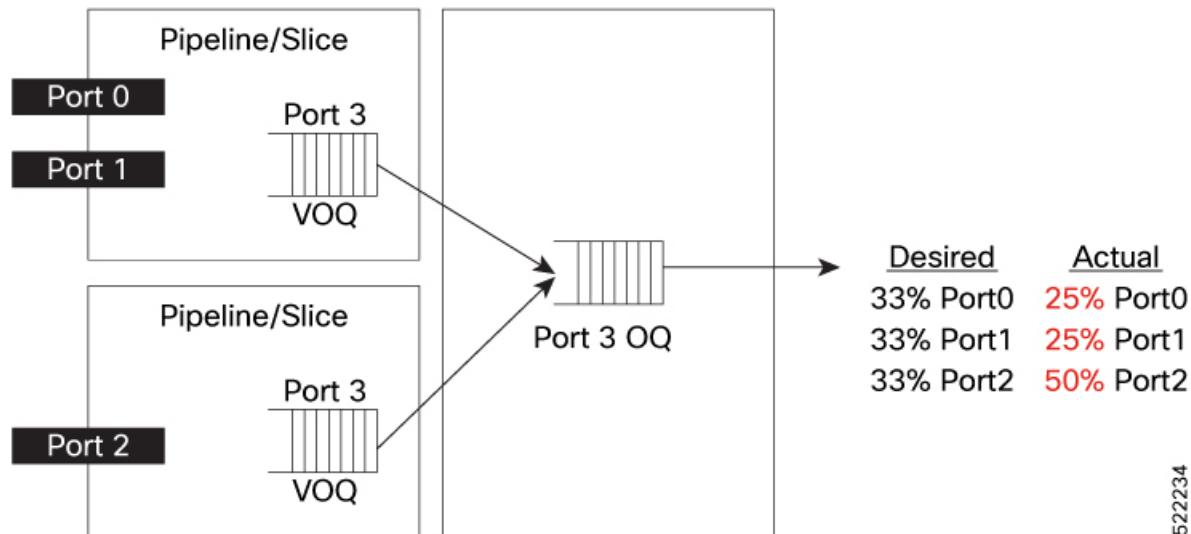
Fair VOQ: Why

Per default behavior, every network slice of an NPU is assigned a set of 4 or 8 Virtual Output Queues (VOQ) per destination port. With such an assignment, it's challenging to ensure that the right amount of buffering is available via VOQs. With this configuration, the ingress traffic from various source ports on a slice (or pipeline) on an NPU destined to a destination port is assigned to a VOQ per slice. In other words, multiple source ports sending traffic to the same destination port use the same VOQ. However, when sending traffic to different destination ports, traffic is enqueued to different VOQs. This means that the traffic isn't distributed equitably because each slice doesn't get its fair share of the output queue bandwidth. In a scenario where one slice has two ports and another slice has only one port, the bandwidth drops for the ports sharing a slice, even though the two ports handle more traffic than the single port.

Consider the following example where two 100G ports—port-0 and port-1—that belong to the same slice (slice-0) are sending traffic to port-3 on the output queue (OQ). You have a 100G port on another slice (slice-1) on the same NPU that is also scheduled to send traffic to port-3. The ingress VOQ is shared between the two

ports in slice-0, whereas the ingress VOQ in slice-1 is available exclusively for port-3. This arrangement results in port-0 and port-1 getting 25% of the buffer traffic, while port-3 receives 50% of the buffer traffic.

Figure 6: Existing behavior : Source ports on slice share one VOQ per destination port



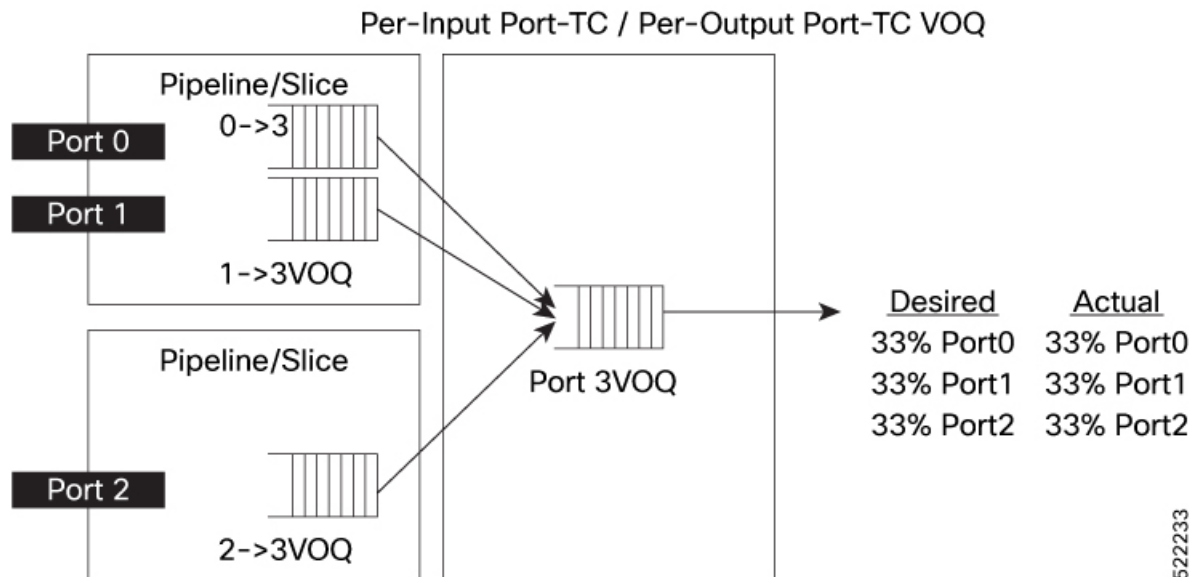
The fair VOQ feature solves this disparity in traffic distribution.

Fair VOQ: How

The fair VOQ feature tackles the default behavior that treats source ports on each NPU slice equally, regardless of the number of active source ports. It does so by redesigning the way bandwidth is allocated from the output queue. Rather than distributing bandwidth at a slice level, fair VOQ distributes bandwidth directly to the source ports. When you configure the command **hw-module profile qos voq-mode** and reload your router, the functionality creates a dedicated VOQ for every source port and destination port pair. This arrangement ensures that bandwidth available at the destination port for a given traffic class is distributed equally to all source ports requesting bandwidth.

Extending the preceding example to understand the fair VOQ functionality, there are now dedicated VOQs for each ingress port that connect to the port on the output queue. Thus, port-0 and port-1 now don't share a VOQ, and port-3 has its VOQ as before, as shown in the following figure. This fair VOQ arrangement results in traffic queued on dedicated queues, thus improving the traffic performance.

Figure 7: Fair VOQ behavior: each source port on slice has one dedicated VOQ per destination port



Fair VOQ Modes and Sharing of Counters

You can configure fair VOQ for 8xVOQ mode (fair-8) and 4xVOQ mode (fair-4) using the following options in the `hw-module profile qos voq-mode` command:

- `hw-module profile qos voq-mode fair-8`
- `hw-module profile qos voq-mode fair-4`

You can also share VOQ statistics counters in both fair VOQ modes, as shown in the following table. (For details on why sharing counters is essential and how to configure counters sharing, see [Sharing of VOQ Statistics Counters](#), on page 45.)

Table 10: Fair VOQ Modes and Sharing Counters

Fair VOQ Mode	Sharing Counters Mode	Important Notes
fair-8	2, 4	<ul style="list-style-type: none"> • Eight VOQs configured per source port and destination pair • Counters are shared by {2, 4} VOQs. • fair-8 mode doesn't support dedicated counter mode (counter mode1, where there's a counter for every queue)

Fair VOQ Mode	Sharing Counters Mode	Important Notes
fair-4	1, 2, 4	<ul style="list-style-type: none"> Four VOQs configured per source port and destination pair Counters are shared by {1, 2, 4} VOQs.

Fair VOQs and Slice (or Normal) VOQs: Key Differences

The following table is a snapshot to outline the key differences between fair VOQs and slice or regular VOQs.

Table 11: Fair VOQs and Normal VOQs

Fair VOQ	Normal VOQ
fair-8 mode: eight VOQs configured per source port and destination pair	8: <ul style="list-style-type: none"> Eight VOQs per destination port per slice These VOQs are shared by all source ports within an NPU slice.
fair-4 mode: four VOQs configured per source port and destination pair	4: <ul style="list-style-type: none"> Four VOQs per destination port per slice These VOQs are shared by all source ports within an NPU slice.

Guidelines and Limitations

- The fair VOQ feature is supported on the Cisco 8202 router (12 QSFP56-DD 400G and 60 QSFP28 100G ports) and Cisco 8201-32FH router (32 QSFP-DD 400G ports).
- The following table details the maximum interfaces (with basic IPv4 configurations and no other scale configuration such as QoS policy, ACL, and subinterface configuration) allowed based on VOQ mode and sharing counter mode.

Table 12: Maximum Interfaces Based on Fair VOQ Mode and Sharing Counter Mode

VOQ Mode	Sharing Counter Mode	Maximum Interfaces for Cisco 8202	Maximum Interfaces for Cisco 8201-32FH
fair-8	1	The router doesn't support this combination.	The router doesn't support this combination.

VOQ Mode	Sharing Counter Mode	Maximum Interfaces for Cisco 8202	Maximum Interfaces for Cisco 8201-32FH
fair-8	2	96 = 60 (100G) + 8x4 + 4 (400G) ==> you can configure only eight 400G interfaces in 4x10G or 4x25G breakout mode.	128 = 32x4 (4x10G or 4x25G - breakout on all 32 ports - 400G)
fair-8	4	108 = 60 + 12 x 4 (breakout on all 12 ports - 400G)	
fair-4	1	96 = 60(100G) + 8x4 + 4 (400G) ==> you can configure only eight 400G interfaces in 4x10G or 4x25G breakout mode.	
fair-4	2	108 = 60 + 12 x4 (breakout on all 12 ports - 400G)	
fair-4	4	108 = 60 + 12 x4 (breakout on all 12 ports - 400G)	



Note We recommend using sharing counter mode 4 in breakout modes and sharing counter mode 2 for nonbreakout modes.

If the counter usage of applications is high, sharing counter mode 2 is not supported.



Note Breakout mode isn't supported on 100G interfaces.

- Ensure that you reload the router for the configuration to take effect.
- Layer 2 traffic isn't supported in **fair-voq** mode (**fair-4** and **fair-8**).
- Subinterface queuing isn't supported. (This applies to bundle sub-interfaces as well). This means that you can't attach egress service-policies that require dedicated VOQs. However, egress marking is supported for subinterfaces.
- ERSPAN mirroring traffic is not supported in **fair-voq** mode.
- **hw-module profile stats voqs-sharing-counters 1** isn't supported in **fair-8** mode. Ensure that you configure **hw-module profile voq sharing-counters 2** or **hw-module profile voq sharing-counters 4**

along with **hw-module profile qos voq-mode fair-4** or **hw-module profile qos voq-mode fair-8** before reloading the router.

- Breakout is supported only on 400G interfaces in **fair-voq** mode (both **fair-4** and **fair-8**) on the Cisco 8202 router.
- **src-interface** and **src-slice** keywords in **show controller npu stats** are visible only when you configure the VOQ mode to either **fair-8** or **fair-4**.

Configure Fair VOQ

To configure fair VOQ:

1. Configure sharing of VOQ statistics counters. This example configures 2 counters.



Note Configuring **fair-8** mode without counter sharing may cause configuration failure or other unexpected behavior.

2. Configure fair VOQ mode. This example shows how to configure **fair-8** mode.
3. Restart the router for the configuration to take effect.
4. You have successfully enabled the fair VOQ feature to ensure equitable traffic distribution between every source port and destination port pair.

```
/*Configure sharing of VOQ statistics counters; we're configuring 2 counters per queue*/
Router(config)#hw-module profile stats ?
  voqs-sharing-counters  Configure number of voqs (1, 2, 4) sharing counters
Router(config)#hw-module profile stats voqs-sharing-counters ?
  1  Counter for each queue
  2  2 Queues share counters
  4  4 Queues share counters
Router(config)#hw-module profile stats voqs-sharing-counters 2

/*Configure fair-voq mode; we're configuring fair-8 VOQ mode here*/
Router#config
Router(config)#hw-module profile qos voq-mode fair-8
Router(config)#commit
Router#reload location all
```

Running Configuration

```
hw-module profile stats voqs-sharing-counters 2
!
hw-module profile qos voq-mode fair-8
!
```

Verification

Run the **show controller npu stats voq ingress interface <> instance <> location <>** command to verify the fair VOQ configuration.

```
Router#show controllers npu stats voq ingress interface hundredGigE 0/0/0/20 instance 0
location 0/RP0/CPU0
```

```
Interface Name      =  Hu0/0/0/20
```

```

Interface Handle      =      f000118
Location              =      0/RP0/CPU0
Asic Instance         =      0
Port Speed(kbps)     =      100000000
Local Port           =      local
Src Interface Name    =      ALL
VOQ Mode            =      Fair-8
Shared Counter Mode  =      2
  ReceivedPkts      ReceivedBytes  DroppedPkts  DroppedBytes
-----
TC_{0,1} = 11110      1422080          0            0
TC_{2,3} = 0           0                0            0
TC_{4,5} = 0           0                0            0
TC_{6,7} = 0           0                0            0
RP/0/RP0/CPU0:ios#

```

Associated Commands

[hw-module profile qos voq-mode](#)

Modular QoS Congestion Avoidance

Congestion avoidance techniques monitor traffic flow in an effort to anticipate and avoid congestion at common network bottlenecks. Avoidance techniques are implemented before congestion occurs as compared with congestion management techniques that control congestion after it has occurred.

Congestion avoidance is achieved through packet dropping. The router supports these QoS congestion avoidance techniques:

- [Tail Drop and the FIFO Queue, on page 54](#)
- [Random Early Detection and TCP, on page 56](#)

Tail Drop and the FIFO Queue

Tail drop is a congestion avoidance technique that drops packets when an output queue is full until congestion is eliminated. Tail drop treats all traffic flow equally and does not differentiate between classes of service. It manages the packets placed into a first-in, first-out (FIFO) queue, and forwarded at a rate determined by the available underlying link bandwidth.

Configure Tail Drop

Packets satisfying the match criteria for a class accumulate in the queue reserved for the class until they are serviced. The **queue-limit** command is used to define the maximum threshold for a class. When the maximum threshold is reached, the enqueued packets to the class queue result in tail drop (packet drop).

Restrictions

- When configuring the **queue-limit** command, you must configure one of the following commands: **priority**, **shape average**, or **bandwidth remaining**, except for the default class.

Configuration Example

You have to accomplish the following to complete the tail drop configuration:

1. Creating (or modifying) a policy map that can be attached to one or more interfaces to specify a service policy
2. Associating the traffic class with the traffic policy
3. Specifying the maximum limit the queue can hold for a class policy configured in a policy map.
4. Specifying priority to a class of traffic belonging to a policy map.
5. (Optional) Specifying the bandwidth allocated for a class belonging to a policy map or specifying how to allocate leftover bandwidth to various classes.
6. Attaching a policy map to an output interface to be used as the service policy for that interface.

```
Router# configure
Router(config)# policy-map test-qlimit-1
Router(config-pmap)# class qos-1
Router(config-pmap-c)# queue-limit 100 us
Router(config-pmap-c)# priority level 7
Router(config-pmap-c)# exit
Router(config-pmap)# exit

Router(config)# interface HundredGigE 0/6/0/18
Router(config-if)# service-policy output test-qlimit-1
Router(config-if)# commit
```

Running Configuration

```
policy-map test-qlimit-1
  class qos-1
    queue-limit 100 us
    priority level 7
  !
  class class-default
  !
end-policy-map
!
```

Verification

```
Router# show qos int hundredGigE 0/6/0/18 output
```

```
NOTE:- Configured values are displayed within parentheses
Interface HundredGigE0/6/0/18 ifh 0x3000220 -- output policy
NPU Id: 3
Total number of classes: 2
Interface Bandwidth: 100000000 kbps
VOQ Base: 11176
VOQ Stats Handle: 0x88550ea0
Accounting Type: Layer1 (Include Layer 1 encapsulation and above)
-----
Level1 Class (HP7) = qos-1
Egressq Queue ID = 11177 (HP7 queue)
```

```
TailDrop Threshold           = 1253376 bytes / 100 us (100 us)
WRED not configured for this class

Level1 Class                 = class-default
Egressq Queue ID            = 11176 (Default LP queue)
Queue Max. BW.              = 101803495 kbps (default)
Queue Min. BW.              = 0 kbps (default)
Inverse Weight / Weight     = 1 (BWR not configured)
TailDrop Threshold          = 1253376 bytes / 10 ms (default)
WRED not configured for this class
```

Related Topics

- [Tail Drop and the FIFO Queue, on page 54](#)

Random Early Detection and TCP

The Random Early Detection (RED) congestion avoidance technique takes advantage of the congestion control mechanism of TCP. By randomly dropping packets prior to periods of high congestion, RED tells the packet source to decrease its transmission rate. Assuming the packet source is using TCP, it decreases its transmission rate until all packets reach their destination, indicating that the congestion is cleared. You can use RED as a way to cause TCP to slow transmission of packets. TCP not only pauses, but it also restarts quickly and adapts its transmission rate to the rate that the network can support.

RED distributes losses in time and maintains normally low queue depth while absorbing traffic bursts. It achieves this by taking action on the average queue size, and not the instantaneous queue size. When enabled on an interface, RED begins dropping packets when congestion occurs at a rate you select during configuration.

Configure Random Early Detection

The **random-detect** command with the minimum threshold and maximum threshold keywords must be used to enable random early detection (RED).

Guidelines

- If you configure the **random-detect <min threshold> <max threshold>** command on any class including class-default, configure one of the following commands: **shape average** or **bandwidth remaining**.
- If you configure a queue-limit that is lesser than the minimum supported value, the configured value automatically adjusts to the supported minimum value.

While configuring **random-detect**, if you set the **<min threshold>** and **<max-threshold>** values lesser than the minimum supported threshold value:

- The **<min threshold>** value automatically adjusts to the minimum supported value.
- The **<max-threshold>** value doesn't autoadjust to a value above the minimum supported threshold value. This results in a failed **random-detect** configuration. To prevent this error, configure the **<max-threshold>** value such that it exceeds the **<min threshold>** value that your system supports.

Configuration Example

Accomplish the following to complete the random early detection configuration:

1. Creating (or modifying) a policy map that can be attached to one or more interfaces to specify a service policy
2. Associating the traffic class with the traffic policy
3. Enabling RED with minimum and maximum thresholds.
4. Configure **one** of the following:
 - Specifying how to allocate leftover bandwidth to various classes. **OR**
 - Shaping traffic to the specified bit rate or a percentage of the available bandwidth.
5. Attaching a policy map to an output interface to be used as the service policy for that interface.

```
Router# configure
Router(config)# policy-map red-abs-policy
Router(config-pmap)# class qos-1
Router(config-pmap-c)# random-detect <min threshold> <max threshold>
Router(config-pmap-c)# shape average percent 10
Router(config-pmap-c)# end-policy-map
Router(config)# commit
Router(config)# interface HundredGigE0/0/0/12
Router(config-if)# service-policy output red-abs-policy
Router(config-if)# commit
```

Running Configuration

```
policy-map red-abs-policy
class tc7
  priority level 1
  queue-limit 75 mbytes
!
class tc6
  priority level 2
  queue-limit 75 mbytes
!
class tc5
  shape average 10 gbps
  queue-limit 75 mbytes
!
class tc4
  shape average 10 gbps
  queue-limit 75 mbytes
!
class tc3
  shape average 10 gbps
  queue-limit 75 mbytes
!
class tc2
  shape average 10 gbps
  queue-limit 75 mbytes
!
class tc1
  shape average 10 gbps
```

```

random-detect ecn
random-detect 100 mbytes 200 mbytes
!
class class-default
shape average 10 gbps
random-detect 100 mbytes 200 mbytes
!
end-policy-map
!

interface HundredGigE0/0/0/12
service-policy output red-abs-policy
shutdown
!
```

Verification

Router# **show qos int hundredGigE 0/6/0/18 output**

```

NOTE:- Configured values are displayed within parentheses
Interface HundredGigE0/0/0/12 ifh 0x3000220 -- output policy
NPU Id:                               3
Total number of classes:               2
Interface Bandwidth:                   100000000 kbps
VOQ Base:                               11176
VOQ Stats Handle:                       0x88550ea0
Accounting Type:                        Layer1 (Include Layer 1 encapsulation and above)
-----
Level1 Class                           = qos-1
Egressq Queue ID                       = 11177 (LP queue)
Queue Max. BW.                          = 10082461 kbps (10 %)
Queue Min. BW.                          = 0 kbps (default)
Inverse Weight / Weight                  = 1 (BWR not configured)
Guaranteed service rate                  = 10000000 kbps
TailDrop Threshold                       = 12517376 bytes / 10 ms (default)

Default RED profile
RED Min. Threshold                       = 12517376 bytes (10 ms)
RED Max. Threshold                       = 12517376 bytes (10 ms)

Level1 Class                             = class-default
Egressq Queue ID                         = 11176 (Default LP queue)
Queue Max. BW.                           = 101803495 kbps (default)
Queue Min. BW.                           = 0 kbps (default)
Inverse Weight / Weight                   = 1 (BWR not configured)
Guaranteed service rate                   = 50000000 kbps
TailDrop Threshold                       = 62652416 bytes / 10 ms (default)
WRED not configured for this class
```

Related Topics

- [Random Early Detection and TCP, on page 56](#)

Explicit Congestion Notification

Random Early Detection (RED) is implemented at the core routers of a network. Edge routers assign IP precedences to packets, as the packets enter the network. With RED, core routers then use these precedences

to determine how to treat different types of traffic. RED provides a single threshold and weights per traffic class or queue for different IP precedences.

ECN is an extension to RED. ECN marks packets instead of dropping them when the average queue length exceeds a specific threshold value. When configured, ECN helps routers and end hosts to understand that the network is congested and slow down sending packets. However, if the queue length is above the maximum threshold for the extended memory, packets are dropped. This is the identical treatment that a packet receives when RED is enabled without ECN configured on the router.

RFC 3168, *The Addition of Explicit Congestion Notification (ECN) to IP*, states that with the addition of active queue management (for example, RED) to the Internet infrastructure, routers are no longer limited to packet loss as an indication of congestion.



Note You cannot use this feature when you have set qos-group or mpls experimental along with a traffic class in the ingress policy.

Implementing ECN

Implementing ECN requires an ECN-specific field that has two bits—the ECN-capable Transport (ECT) bit and the CE (Congestion Experienced) bit—in the IP header. The ECT bit and the CE bit can be used to make four code points of 00 to 11. The first number is the ECT bit and the second number is the CE bit.

Table 13: ECN Bit Setting

ECT Bit	CE Bit	Combination Indicates
0	0	Not-ECN-capable.
0	1	Endpoints of the transport protocol are ECN-capable.
1	0	Endpoints of the transport protocol are ECN-capable.
1	1	Congestion experienced.

The ECN field combination 00 indicates that a packet is not using ECN. The code points 01 and 10—Called ECT(1) and ECT(0), respectively—are set by the data sender to indicate that the endpoints of the transport protocol are ECN-capable. Routers treat these two code points identically. Data senders can use either one or both of these two combinations. The ECN field combination 11 indicates congestion to the endpoints. Packets arriving a full queue of a router will be dropped.

Packet Handling When ECN Is Enabled

When ECN is enabled, all packets between <min_threshold> and <max tail drop threshold> are marked with ECN. Three different scenarios arise if the queue length is between the minimum threshold and the maximum threshold:

- If the ECN field on the packet indicates that the endpoints are ECN-capable (that is, the ECT bit is set to 1 and the CE bit is set to 0, or the ECT bit is set to 0 and the CE bit is set to 1)—and the RED algorithm determines that the packet should have been dropped based on the drop probability—the ECT and CE

bits for the packet are changed to 1, and the packet is transmitted. This happens because ECN is enabled and the packet gets marked instead of dropped.

- If the ECN field on the packet indicates that neither endpoint is ECN-capable (that is, the ECT bit is set to 0 and the CE bit is set to 0), the packet is transmitted. If, however, the max tail drop threshold is exceeded, the packet is dropped. This is the identical treatment that a packet receives when RED is enabled without ECN configured on the router.
- If the ECN field on the packet indicates that the network is experiencing congestion (that is, both the ECT bit and the CE bit are set to 1), the packet is transmitted. No further marking is required.

Configuration Example

```
Router# configure
Router(config)# policy-map policy1
Router(config-pmap)# class class1
Router(config-pmap-c)# bandwidth percent 50
Router(config-pmap-c)# random-detect 1000 packets 2000 packets
Router(config-pmap-c)# random-detect ecn
Router(config-pmap-c)# exit
Router(config-pmap)# exit
Router(config)# commit
```

Verification

Use the **show policy-map interface** to verify the configuration.

```
Router# show policy-map int hu 0/0/0/35 output
TenGigE0/0/0/6 output: pm-out-queue

HundredGigE0/0/0/35 output: egress_qosgrp_ecn

Class tc7
  Classification statistics          (packets/bytes)      (rate - kbps)
  Matched                          : 195987503/200691203072    0
  Transmitted                       : 188830570/193362503680    0
  Total Dropped                     : 7156933/7328699392       0
  Queueing statistics
  Queue ID                          : 18183
  Taildropped(packets/bytes)        : 7156933/7328699392

  WRED profile for
  RED Transmitted (packets/bytes)    : N/A
  RED random drops (packets/bytes)   : N/A
  RED maxthreshold drops (packets/bytes) : N/A
  RED ecn marked & transmitted(packets/bytes): 188696802/193225525248

Class tc6
  Classification statistics          (packets/bytes)      (rate - kbps)
  Matched                          : 666803815/133360763000    0
  Transmitted                       : 642172362/128434472400    0
  Total Dropped                     : 24631453/4926290600     0
  Queueing statistics
  Queue ID                          : 18182
  Taildropped(packets/bytes)        : 24631453/4926290600

  WRED profile for
  RED Transmitted (packets/bytes)    : N/A
  RED random drops (packets/bytes)   : N/A
  RED maxthreshold drops (packets/bytes) : N/A
  RED ecn marked & transmitted(packets/bytes): 641807908/128361581600
```

```
Class tc5
Classification statistics          (packets/bytes)      (rate - kbps)
  Matched                        :      413636363/82727272600      6138
  Transmitted                     :      398742312/79748462400      5903
  Total Dropped                   :      14894051/2978810200       235
Queueing statistics
  Queue ID                       :      18181
  Taildropped (packets/bytes)     :      14894051/2978810200

WRED profile for
RED Transmitted (packets/bytes)   : N/A
RED random drops (packets/bytes)  : N/A
RED maxthreshold drops (packets/bytes) : N/A
RED ecn marked & transmitted (packets/bytes): 398377929/79675585800
```



Note The **RED ecn marked & transmitted(packets/bytes)** row displays the statistics for ECN marked packets. To begin with, it displays *0/0*.



CHAPTER 6

Configure Priority Flow Control

- [Priority Flow Control Overview](#), on page 63
- [View Priority Flow Control TX Frames Per Traffic Class](#) , on page 72
- [Configurable ECN Threshold and Maximum Marking Probability Values](#), on page 74
- [Priority Flow Control Watchdog Overview](#), on page 79
- [Global Statistics Counters for Priority Flow Control and Priority Flow Control Watchdog](#), on page 82
- [PFC Watchdog Statistics to Monitor Traffic Drops](#) , on page 85
- [Detect High Bandwidth Memory Congestion](#), on page 87

Priority Flow Control Overview

Table 14: Feature History Table

Feature Name	Release Information	Feature Description
Priority Flow Control on Cisco 8808 and Cisco 8812 Modular Chassis Line Cards	Release 7.5.3	<p>Priority Flow Control is now supported on the following line card in the buffer-internal mode:</p> <ul style="list-style-type: none"> • 88-LC0-34H14FH <p>The feature is supported in the buffer-internal and buffer-extended modes on:</p> <ul style="list-style-type: none"> • 88-LC0-36FH <p>Apart from the buffer-external mode, support for this feature now extends to the buffer-internal mode on the following line cards:</p> <ul style="list-style-type: none"> • 88-LC0-36FH-M • 8800-LC-48H

Feature Name	Release Information	Feature Description
Shortlink Priority Flow Control	Release 7.3.3	This feature and the hw-module profile priority-flow-control command are supported on 88-LC0-36FH line card.
Priority Flow Control Support on Cisco 8800 36x400 GbE QSFP56-DD Line Cards (88-LC0-36FH-M)	Release 7.3.15	<p>This feature and the hw-module profile priority-flow-control command are supported on 88-LC0-36FH-M and 8800-LC-48H line cards.</p> <p>All previous functionalities and benefits of this feature are available on these line cards. However, the buffer-internal mode is not supported.</p> <p>In addition, to use the buffer-extended mode on these line cards, you are required to configure the performance capacity or headroom values. This configuration requirement ensures that you can better provision and balance workloads to achieve lossless behavior, which in turn ensures efficient use of bandwidth and resources.</p>
Priority Flow Control	Release 7.3.1	This feature and the hw-module profile priority-flow-control command are not supported.

Priority-based Flow Control (IEEE 802.1Qbb), which is also referred to as Class-based Flow Control (CBFC) or Per Priority Pause (PPP), is a mechanism that prevents frame loss that is due to congestion. PFC is similar to 802.x Flow Control (pause frames) or link-level flow control (LFC). However, PFC functions on a per class-of-service (CoS) basis.

During congestion, PFC sends a pause frame to indicate the CoS value to pause. A PFC pause frame contains a 2-octet timer value for each CoS that indicates the length of time to pause the traffic. The unit of time for the timer is specified in pause quanta. A quanta is the time required for transmitting 512 bits at the speed of the port. The range is from 0 through 65535 quanta.

PFC asks the peer to stop sending frames of a particular CoS value by sending a pause frame to a well-known multicast address. This pause frame is a one-hop frame and isn't forwarded when received by the peer. When the congestion mitigates, the router stops sending the PFC frames to the upstream node.

You can configure PFC for each line card using the **hw-module profile priority-flow-control** command in one of two modes:

- buffer-internal
- buffer-extended



Note PFC threshold configurations are deprecated in `pause` command. Use the `hw-module profile priority-flow-control` command to configure PFC threshold configurations.

Related Topics

- [Configure Priority Flow Control, on page 67](#)
- [Priority Flow Control Watchdog Overview, on page 79](#)

buffer-internal mode

Use this mode if PFC-enabled devices aren't more than 1 km apart.

You can set values for pause-threshold, headroom (both related to PFC), and ECN for the traffic class using the `hw-module profile priority-flow-control` command in this mode. The buffer-internal configuration applies to all ports that the line card hosts, which mean that you can configure a set of these values per line card.

The existing queue limit and ECN configuration in the queuing policy attached to the interface has no impact in this mode.

The effective queue limit for this mode = pause-threshold + headroom (in bytes)

Restrictions and Guidelines

The following restrictions and guidelines apply while configuring the PFC threshold values using the buffer-internal mode.

- The PFC feature isn't supported on fixed chassis systems.
- Ensure that there's no breakout configured on a chassis that has the PFC configured. Configuring PFC and breakout on the same chassis may lead to unexpected behavior, including traffic loss.
- The feature isn't supported on bundle and non-bundle sub-interface queues.
- The feature is supported on 40GbE, 100 GbE, and 400 GbE interfaces.
- The feature isn't supported in the 4xVOQ queuing mode.
- The feature isn't supported when sharing of VOQ counters is configured.

buffer-extended mode

Use this mode for PFC-enabled devices with long-haul connections.

You can set the value for pause-threshold using the `hw-module profile priority-flow-control` command in this mode. You must, however, configure the queuing policy attached to the interface to set the ECN and queuing limits. The buffer-extended configuration applies to all ports that the line card hosts, which mean that you can configure a set of these values per line card.

Configuration Guidelines

- Important points while configuring the buffer-extended mode on 88-LC0-36FH-M line cards:
 - Apart from pause-threshold, you must also configure values for headroom.
 - The headroom value range is from 4 through 75000.
 - Specify pause-threshold and headroom values in units of kilobytes (KB) or megabytes (MB).
- Important points while configuring the buffer-extended mode on 8800-LC-48H line cards:
 - Configure values only for **pause-threshold**. Don't configure headroom values.
 - Configure **pause-threshold** in units of milliseconds (ms) or microseconds.
 - Don't use units of kilobytes (KB) or megabytes (MB) units, even though the CLI displays them as options. Only use units of milliseconds (ms) or microseconds.

(Also see [Configure Priority Flow Control](#), on page 67)

Important Considerations

- If you configure PFC values in the buffer-internal mode, then the ECN value for the line card is derived from the buffer-internal configuration. If you configure PFC values in the buffer-extended mode, then the ECN value is derived from the policy map. (For details on the ECN feature, see [Explicit Congestion Notification](#), on page 58.)
- The buffer-internal and buffer-extended modes can't coexist on the same line card.
- If you add or remove traffic-class actions on a line card, you must reload the line card.
- When using the buffer-internal mode, you can change values of the following parameters without having to reload the line card. However, if you add a new traffic class and configure these values for the first time on that traffic class, you must reload the line card for the values to come into effect.
 - pause-threshold
 - headroom
 - ECN
- If you add or remove ECN configuration using the **hw-module profile priority-flow-control** command, you must reload the line card for the ECN changes to take effect.
- The PFC threshold value ranges for the buffer-internal mode are as follows.

Threshold	Configured (bytes)
pause (min)	307200
pause (max)	422400
headroom (min)	345600
headroom (max)	537600

Threshold	Configured (bytes)
ecn (min)	153600
ecn (max)	403200

- For a traffic-class, the ECN value must always be lesser than the configured pause-threshold value.
- The combined configured values for pause-threshold and headroom must not exceed 844800 bytes. Else, the configuration is rejected.
- The pause-threshold value range for buffer-extended mode is from 2 milliseconds (ms) through 25 ms and from 2000 microseconds through 25000 microseconds.

Hardware Support for Priority Flow Control

The table lists the PIDs that support PFC per release and the PFC mode in which the support is available.

Table 15: PFC Hardware Support Matrix

Release	PID	PFC Mode
Release 7.5.3	88-LC0-36FH	buffer-extended and buffer-internal
	88-LC0-36FH-M	buffer-extended and buffer-internal
	8800-LC-48H	buffer-extended and buffer-internal
	88-LC0-34H14FH	buffer-internal
Release 7.3.15	<ul style="list-style-type: none"> • 88-LC0-36FH-M • 88-LC0-36FH 	buffer-extended
Release 7.0.11	8800-LC-48H	buffer-internal

Configure Priority Flow Control

You can configure PFC to enable the no-drop behavior for the CoS as defined by the active network QoS policy.



Note The system enables shortlink PFC by default when you enable PFC.

Configuration Example

You must accomplish the following to complete the PFC configuration:

1. Enable PFC at the interface level.
2. Configure ingress classification policy.

3. Attach the PFC policy to the interface.
4. Configure PFC threshold values using either the buffer-internal or buffer-extended mode.

```

Router# configure
Router(config)# priority-flow-control mode on
/*Configure ingress classification policy*/
Router(config)# class-map match-any prec7
Router(config-cmap)# match precedence
Router(config)# class-map match-any tc7
/*Ingress policy attach*/
Router(config-if)# service-policy input QOS_marking
/*Egress policy attach*/
Router(config-if)# service-policy output qos_queuing
Router(config-pmap-c)# exit
Router(config-pmap)# exit
Router(config)#show controllers npu priority-flow-control location <loc>

```

Running Configuration

```

*Interface Level*
interface HundredGigE0/0/0/0
    priority-flow-control mode on

*Ingress:*
class-map match-any prec7

    match precedence 7

end-class-map

!

class-map match-any prec6

    match precedence 6

end-class-map

!

class-map match-any prec5

    match precedence 5

end-class-map

!

class-map match-any prec4

    match precedence 4

end-class-map

!

class-map match-any prec3
    match precedence 3
end-class-map
!
class-map match-any prec2

```

```
match precedence 2
end-class-map
!
class-map match-any prec1
match precedence 1
end-class-map
!
!
policy-map QOS_MARKING
class prec7
    set traffic-class 7
    set qos-group 7
!
class prec6
    set traffic-class 6
    set qos-group 6
!
class prec5
    set traffic-class 5
    set qos-group 5
!
class prec4
    set traffic-class 4
    set qos-group 4
!
class prec3
    set traffic-class 3
    set qos-group 3
!
class prec2
    set traffic-class 2
    set qos-group 2
!
class prec1
    set traffic-class 1
    set qos-group 1
!
class class-default
    set traffic-class 0
    set qos-group 0
!

*Egress:*
class-map match-any tc7
match traffic-class 7
end-class-map
!
class-map match-any tc6
match traffic-class 6
end-class-map
!
class-map match-any tc5
match traffic-class 5
end-class-map

!

class-map match-any tc4

match traffic-class 4

end-class-map

!
```

```

class-map match-any tc3
  match traffic-class 3
end-class-map

!

class-map match-any tc2
match traffic-class 2
end-class-map
!
class-map match-any tc1
match traffic-class 1
end-class-map
!
policy-map QOS_QUEUING
class tc7
  priority level 1
  shape average percent 10
!
class tc6
  bandwidth remaining ratio 1
  queue-limit 100 ms
!
class tc5
  bandwidth remaining ratio 20
  queue-limit 100 ms
!
class tc4
  bandwidth remaining ratio 20
  random-detect ecn
  random-detect 6144 bytes 100 mbytes
!
class tc3
  bandwidth remaining ratio 20
  random-detect ecn
  random-detect 6144 bytes 100 mbytes
!
class tc2
  bandwidth remaining ratio 5
  queue-limit 100 ms
!
class tc1
  bandwidth remaining ratio 5
  queue-limit 100 ms
!
class class-default
  bandwidth remaining ratio 20
  queue-limit 100 ms
!
[buffer-extended]

hw-module profile priority-flow-control location 0/0/CPU0
  buffer-extended traffic-class 3 pause-threshold 10 ms
  buffer-extended traffic-class 4 pause-threshold 10 ms
!

[buffer-internal]

hw-module profile priority-flow-control location 0/1/CPU0
  buffer-internal traffic-class 3 pause-threshold 403200 bytes headroom 441600 bytes ecn
  224640 bytes

```



```
buffer-internal traffic-class 4 pause-threshold 403200 bytes headroom 441600 bytes ecn
224640 bytes
```

Verification

```
Router#sh controllers hundredGigE0/0/0/22 priority-flow-control
Priority flow control information for interface HundredGigE0/0/0/22:
```

```
Priority Flow Control:
Total Rx PFC Frames : 0
Total Tx PFC Frames : 313866
Rx Data Frames Dropped: 0
CoS Status Rx Frames
```

```
--- -----
0 on 0
1 on 0
2 on 0
3 on 0
4 on 0
5 on 0
6 on 0
7 on 0
```

```
/*[buffer-internal]*/
```

```
Router#show controllers hundredGigE 0/9/0/24 priority-flow-control
```

```
Priority flow control information for interface HundredGigE0/9/0/24:
```

```
Priority Flow Control:
Total Rx PFC Frames : 0
Total Tx PFC Frames : 313866
Rx Data Frames Dropped: 0
CoS Status Rx Frames
```

```
--- -----
0 on 0
1 on 0
2 on 0
3 on 0
4 on 0
5 on 0
6 on 0
7 on 0
...
```

```
/*[buffer-internal, tc3 & tc4 configured. TC4 doesn't have ECN]*/
```

```
Router#show controllers npu priority-flow-control location <loc>
```

```
Location Id: 0/1/CPU0
PFC: Enabled
PFC-Mode: buffer-internal
TC Pause Headroom ECN
-----
3 86800 bytes 120000 bytes 76800 bytes
4 86800 bytes 120000 bytes Not-configured
```

```
/*[buffer-extended PFC, tc3 & tc4 configured]*/
```

```
Router#show controllers npu priority-flow-control location <loc>
```

```
Location Id: 0/1/CPU0
PFC: Enabled
PFC-Mode: buffer-extended
TC Pause
-----
```

```

3      5000 us
4      10000 us

/*[No PFC]*/

Router#show controllers npu priority-flow-control location <loc>
Location Id:                0/1/CPU0
PFC:                        Disabled

```

Related Topics

- [Priority Flow Control Overview, on page 63](#)

Related Commands hw-module profile priority-flow-control location

View Priority Flow Control TX Frames Per Traffic Class

Table 16: Feature History Table

Feature Name	Release Information	Feature Description
View Priority Flow Control TX Frames Per Traffic Class	Release 7.5.4	<p>You can now view an estimation of the transmission (Tx) of Priority Flow Control (PFC) frames per traffic class, which informs you that the traffic flow has crossed configured PFC pause thresholds. This information allows you to rebalance traffic flows to ensure that network resources are efficiently used. Such verification is possible because we've added counters for PFC Tx frames.</p> <p>This functionality modifies the following to add the PFC Tx frame counters:</p> <ul style="list-style-type: none"> • YANG data model (at Github under the 754 folder): Cisco-IOS-XR-ofa-npu-pfc-oper • CLI: show controllers priority-flow-control statistics

Beginning Cisco IOS XR Release 7.5.4, we have added an additional counter for PFC Tx pause frames per traffic class in the **show controllers priority-flow-control statistics** command. With this counter, you receive timely information that the traffic flow has exceeded the configured PFC pause thresholds and hence that PFC Tx frames are sent out of the PFC-enabled interfaces sourcing such traffic on that network processing core. You can use this information to investigate and analyze why such traffic flows are causing congestion and improve traffic distribution across interfaces to utilize network resources efficiently.

Run the **show controllers priority-flow-control statistics** to view an estimation of the statistics for PFC Tx frames per traffic class (**Tx Frames**):

```
Router#show controllers hundredGigE 0/0/0/4 priority-flow-control statistics
```

Priority flow control information for interface HundredGigE0/0/0/4:
Priority Flow Control:

Total Rx PFC Frames: 0
Total Tx PFC Frames: 4832680
Rx Data Frames Dropped: 1442056 (possible overflow)

CoS	Status	Rx Frames	Tx Frames
0	on	0	0
1	on	0	0
2	on	0	0
3	on	0	2416374
4	on	0	2416306
5	on	0	0
6	on	0	0
7	on	0	0

Guidelines and Limitations for Viewing Priority Flow Control TX Frames Per Traffic Class

This functionality is supported on:

- 8800-LC-36FH-M
- 8800-LC-36FH
- 88-LC0-34H14FH

Configurable ECN Threshold and Maximum Marking Probability Values

Table 17: Feature History Table

Feature Name	Release Information	Feature Description
Configurable ECN Threshold and Maximum Marking Probability Values	Release 7.5.4	<p>While configuring PFC in the buffer-internal mode, you can now optimize congestion notification from the end router to the transmitting router, thus preventing the aggressive throttle of the source traffic. This optimization is possible because we've provided the flexibility to configure the minimum and maximum values for the ECN threshold and the maximum value for marking probability. With these values configured, the probability percentage marking is applied linearly, starting at the ECN minimum threshold till the ECN Max threshold.</p> <p>Earlier releases fixed the maximum ECN marking probability at 100% at the maximum ECN threshold.</p> <p>This functionality adds the following options to the hw-module profile priority-flow-control command:</p> <ul style="list-style-type: none"> • max-threshold • probability-percentage

ECN Threshold and Maximum Marking Probability Values

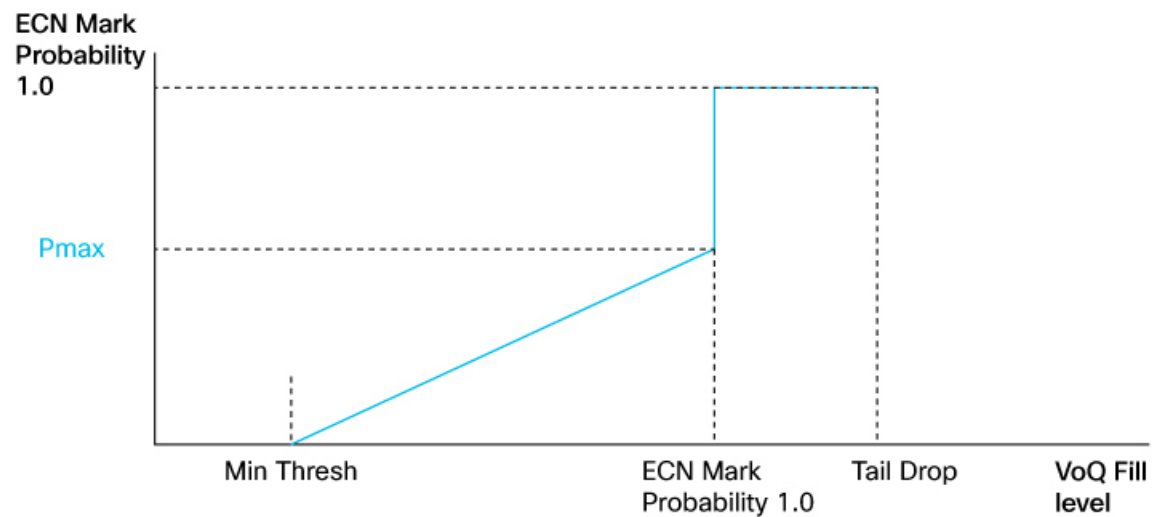
So far, the maximum ECN marking probability was not configurable and was fixed at 100%. You couldn't configure the ECN maximum threshold value either. Such an arrangement of preset marking probabilities and fixed maximum threshold values meant traffic rates began to drop as a function of the queue length. Because of the linear increase in the ECN marking probability—and the consequent congestion signaling from the end host to the transmitting host—traffic rates could begin slowing down even if your link had the necessary bandwidth.

With this feature, you can configure the minimum and maximum ECN threshold values and, depending on your link capability and traffic requirements, configure a marking probability such that the number of packets being marked between the minimum and maximum thresholds is reduced.

The figure depicts the feature implementation visually, where:

- Pmax denotes the configurable maximum marking probability.
- The ECN mark probability is 1 from the maximum threshold to the tail drop threshold.
- Mark probability increases linearly from 0 at the minimum threshold to Pmax at the maximum threshold.

Figure 8: ECN Mark Probability vs. Queue Length (VOQ Fill Levels)



You configure the **max-threshold** and **probability-percentage** options for this feature in the **hw-module profile priority-flow-control** command. We have provided you the flexibility to choose one of these configuration options:

- Default **max-threshold** and **probability-percentage** values.
- Configurable **max-threshold** and **probability-percentage** values.
- PFC in the buffer-internal mode without the new options, as you did before Release 7.5.4.

See [Configure ECN Threshold and Maximum Marking Probability Values, on page 77](#) for details.

Benefits of Configurable ECN Threshold and Maximum Marking Probability Values

A lower mark probability allows the end host to better identify the aggressor flow in the queue and only throttle it without impacting other well-behaved flows. This in turn leads to optimized and efficient deployment and usage of bandwidth.

ECN Threshold and Maximum Marking Probability Values: FAQs

- What if I set the ECN marking probability percentage to 100%?

ECN Marking is linear once queue length exceeds the ECN minimum value. The tail drop comes into play when the average queue size reaches the marking probability percentage of 100%.

- **How does traffic behave beyond my configured value of maximum ECN marking probability?**

Let us say you've set the maximum ECN marking probability to 5%. Any more increase in the average queue length beyond 5% shifts the marking probability to 100%, and [Tail Drop and the FIFO Queue, on page 54](#) comes into effect.

Guidelines and Limitations

- This functionality is available only when you configure PFC in the buffer-internal mode.
- If you configure PFC values in the buffer-internal mode, the ECN value for the line card is derived from the buffer-internal configuration, unlike in the buffer-extended mode where the ECN value is derived from the policy map.
- Line Cards that support this functionality:
 - 88-LC0-36FH
 - 88-LC0-36FH-M
- The following interface types support this functionality:
 - Physical interfaces
 - Bundle interfaces
 - Subinterfaces
 - Bundle subinterfaces
- This functionality is supported for all interface speeds.
- If you have a policy map with one or multiple classes with maximum ECN marking probability enabled, you can:
 - Apply the map to any of the supported interface types.
 - Remove the map from any of the supported interface types.
 - Modify the map while you're attaching it to multiple interfaces.
- The probability percentage option is supported only with **random-detect ecn** configured in the same class. Else, the policy is rejected when applied on an interface. **(Not sure this is applicable. Please confirm.)**
- Ensure that the probability percentage configuration is the same for all traffic classes because this is a device-level configuration.

Table 18: Useful Tips

If you...	...you must...
<p>Want to switch from the default configuration mode to the custom configuration mode</p> <p>OR</p> <p>Want to switch from the custom configuration mode to the default configuration</p>	<ol style="list-style-type: none"> 1. Use the no form of the hw-module profile priority-flow-control command to remove the existing configuration. 2. Configure the new mode and settings using the hw-module profile priority-flow-control command. 3. Reload the line card
<p>Configured PFC in buffer-internal mode without configuring the max-threshold and probability-percentage parameters, but now want to configure them</p>	<ol style="list-style-type: none"> 1. Configure the max-threshold and probability-percentage parameters using the hw-module profile priority-flow-control command. 2. Reload the line card.
<p>Want to change the buffer-internal parameters in the custom mode</p>	<p>Configure the buffer-internal parameters using the hw-module profile priority-flow-control command. You needn't reload the line card.</p>
<p>Want to continue configuring PFC in buffer-internal mode the way you did in releases prior to Release 7.5.4.</p>	<p>Configure the buffer-internal parameters using the hw-module profile priority-flow-control command, but ensure that you only configure values for pause-threshold, headroom, and ecn. If you don't configure values for max-threshold, the router takes the ecn value as the ECN maximum threshold value.</p> <p>See Option 3 under Configure ECN Threshold and Maximum Marking Probability Values, on page 77 for details.</p>

Configure ECN Threshold and Maximum Marking Probability Values

To configure the ECN maximum threshold and maximum probability values:

1. Enable PFC at the interface level.

```
Router(config)#int fourHundredGigE 0/6/0/1
Router(config-if)#priority-flow-control mode on
```

2. [Option 1: Default configuration mode] Configure PFC in buffer-internal mode with pre-defined buffer values enabled along with default **max-threshold** and **probability-percentage** values.

```
Router(config)#hw-module profile priority-flow-control location 0/6/0/1
Router(config-pfc-loc)#buffer-extended traffic-class 3
Router(config-pfc-loc)#buffer-extended traffic-class 4
Router(config-pfc-loc)#commit
```

[Option 2: Custom configuration mode] Configure PFC in buffer-internal mode with custom values for all parameters including **max-threshold** and **probability-percentage**.

```
Router(config)#hw-module profile priority-flow-control location 0/6/0/1
Router(config-pfc-loc)#buffer-internal traffic-class 3 pause-threshold 1574400 bytes
headroom 1651200 bytes ecn 629760 bytes max-threshold 1416960 bytes
probability-percentage 50
Router(config-pfc-loc)#buffer-internal traffic-class 4 pause-threshold 1574400 bytes
headroom 1651200 bytes ecn 629760 bytes max-threshold 1416960 bytes
probability-percentage 50
Router(config-pfc-loc)#commit
```

[Option 3: Without **max-threshold** and **probability-percentage** parameters] Based on your requirements, you can also configure PFC in buffer-internal mode without these parameters.

```
Router(config)#hw-module profile priority-flow-control location 0/6/0/1
Router(config-pfc-loc)#buffer-internal traffic-class 3 pause-threshold 1574400 bytes
headroom 1651200 bytes ecn 629760 bytes
Router(config-pfc-loc)#buffer-internal traffic-class 4 pause-threshold 1574400 bytes
headroom 1651200 bytes ecn 629760 bytes
Router(config-pfc-loc)#commit
```



Note See **Useful Tips** at [ECN Threshold and Maximum Marking Probability Values: FAQs, on page 75](#) to know more about the nuances of switching between configuration modes.

Running Configuration

[Option 1: Default configuration mode]

```
/*Interface Level*/
priority-flow-control mode on
!
/*Line card*/
hw-module profile priority-flow-control location 0/6/0/1
  buffer-extended traffic-class 3
  buffer-extended traffic-class 4
!
```

[Option 2: Custom configuration mode]

```
/*Interface Level*/
priority-flow-control mode on
!
/*Line card*/

hw-module profile priority-flow-control location 0/6/0/1
  buffer-internal traffic-class 3 pause-threshold 1574400 bytes headroom 1651200 bytes ecn
  629760 bytes max-threshold 1416960 bytes probability-percentage 50
  buffer-internal traffic-class 4 pause-threshold 1574400 bytes headroom 1651200 bytes ecn
  629760 bytes max-threshold 1416960 bytes probability-percentage 50
!
```

[Option 3: Without **max-threshold** and **probability-percentage** parameters]

```
/*Interface Level*/
priority-flow-control mode on
!
/*Line card*/
  buffer-internal traffic-class 3 pause-threshold 1574400 bytes headroom 1651200 bytes ecn
  629760 bytes
```



```
buffer-internal traffic-class 4 pause-threshold 1574400 bytes headroom 1651200 bytes ecn
629760 bytes
!
```

Verification

[Option 1: Default configuration mode]

```
Router#show controllers npu priority-flow-control location all
```

```
Location:          0/6/CPU0
PFC:              Enabled
PFC Mode:        buffer-internal
TC   Pause-threshold  Headroom      ECN          ECN-MAX      Prob-per
-----
3    1574400 bytes    1651200 bytes  629760 bytes 1416960 bytes 5
4    1574400 bytes    1651200 bytes  629760 bytes 1416960 bytes 5
```

[Option 2: Custom configuration mode]

```
Router#show controllers npu priority-flow-control location all
```

```
Location:          0/6/CPU0
PFC:              Enabled
PFC Mode:        buffer-internal
TC   Pause-threshold  Headroom      ECN          ECN-MAX      Prob-per
-----
3    1574400 bytes    1651200 bytes  629760 bytes 1416960 bytes 50
4    1574400 bytes    1651200 bytes  629760 bytes 1416960 bytes 50
```

[Option 3: Without **max-threshold** and **probability-percentage** parameters]

```
Router#show controllers npu priority-flow-control location all
```

```
Location:          0/6/CPU0
PFC:              Enabled
PFC Mode:        buffer-internal
TC   Pause-threshold  Headroom      ECN          ECN-MAX      Prob-per
-----
3    1574400 bytes    1651200 bytes  629760 bytes not-configured not-configured
4    1574400 bytes    1651200 bytes  629760 bytes not-configured not-configured
```

Priority Flow Control Watchdog Overview

PFC Watchdog is a mechanism to identify any PFC storms (queue-stuck condition) in the network. It also prevents the PFC from propagating on the network and running in a loop. You can configure a PFC watchdog interval to detect whether packets in a no-drop queue are drained within a specified time period. When the time period is exceeded, all outgoing packets are dropped on interfaces that match the PFC queue that is not being drained.

This requires monitoring PFC receiving on each port and detecting ports seeing an unusual number of sustained pause frames. Once detected, the watchdog module can enforce several actions on such ports, which include generating a syslog message for network management systems, shutting down the queue, and restoring the queue (after the PFC storm stops).

Here's how the PFC Watchdog works:

1. The Watchdog module monitors the PFC-enabled queues to determine the reception of an unusual amount of PFC pause frames in a given interval (Watchdog interval.)
2. Your hardware notifies the Watchdog module when too many PFC frames are received and traffic on the corresponding queues is halted for a time interval.

3. On receiving such notifications, the Watchdog module starts the shutdown timer and moves the queue state to wait-to-shutdown state.
4. At regular intervals during the shutdown interval, the queue is checked for PFC frames and if the traffic in the queue is stuck. If the traffic isn't stuck because the queue didn't receive any PFC frames, the queue moves back to the monitored state.
5. If the traffic is stuck for a longer time and the shutdown-timer expires, the queue switches to a drop state and the PFC Watchdog begins to drop all packets.
6. At regular intervals, the Watchdog checks the queue for PFC frames and whether the traffic in the queue is still stuck. If traffic is stuck in the queue as PFC packets keep arriving, the queue remains in the drop or shutdown state.
7. When the traffic's no longer stuck, the autorestore timer starts. At regular intervals, the module checks if the queue is stuck because of PFC frames.
8. If the queue receives PFC frames during the last autorestore interval, the auto-restore timer is reset upon expiry.
9. If the queue receives no PFC frames during the last autorestore interval, the Watchdog module restores the queue, and traffic resumes.

Related Topics

- [Priority Flow Control Overview, on page 63](#)

Configure a Priority Flow Control Watchdog Interval

You can configure PFC Watchdog parameters (Watchdog interval, shutdown multiplier, auto-restore multiplier) at the global or interface levels. Note that:

- When global Watchdog mode is disabled or off, Watchdog is disabled on all interfaces. This condition is regardless of the interface level Watchdog mode settings.
- When global Watchdog mode is enabled or on, the interface level Watchdog mode configuration settings override the global Watchdog mode values.
- When you configure interface level Watchdog attributes such as interval, shutdown multiplier, and auto-restore multiplier, they override the global Watchdog attributes.



Note Configuring the PFC mode and its policies is a prerequisite for PFC Watchdog.

Configuration Example

You can configure the Watchdog at the global or at the interface level.



Note Watchdog is enabled by default, with system default values of:

Watchdog interval = 100 ms

Shutdown multiplier = 1

Auto-restart multiplier = 10

```
RP/0/RP0/CPU0:ios#show controllers hundredGigE 0/2/0/0 priority-flow-control

Priority flow control information for interface HundredGigE0/2/0/0:

Priority flow control watchdog configuration:
(D) : Default value
U : Unconfigured
-----
Configuration Item          Global  Interface  Effective
-----
PFC watchdog state :      U       U          Enabled(D)
Poll interval :           U       U          100 (D)
Shutdown multiplier :     U       U          1 (D)
Auto-restore multiplier :  U       U          10 (D)
RP/0/RP0/CPU0:ios#config
RP/0/RP0/CPU0:ios(config)#priority-flow-control watchdog mode off
RP/0/RP0/CPU0:ios(config)#commit

RP/0/RP0/CPU0:ios(config)#do show controllers hundredGigE 0/2/0/0 priority-flo$

Priority flow control information for interface HundredGigE0/2/0/0:

Priority flow control watchdog configuration:
(D) : Default value
U : Unconfigured
-----
Configuration Item          Global  Interface  Effective
-----
PFC watchdog state :      Disabled U          Disabled
Poll interval :           U       U          100 (D)
Shutdown multiplier :     U       U          1 (D)
Auto-restore multiplier :  U       U          10 (D)
RP/0/RP0/CPU0:ios(config)#interface hundredGigE 0/2/0/0 priority-flow-control $
RP/0/RP0/CPU0:ios(config)#commit

RP/0/RP0/CPU0:ios(config)#do show controllers hundredGigE 0/2/0/0 priority-flo$

Priority flow control information for interface HundredGigE0/2/0/0:

Priority flow control watchdog configuration:
(D) : Default value
U : Unconfigured
-----
Configuration Item          Global  Interface  Effective
-----
PFC watchdog state :      Disabled Enabled    Disabled
Poll interval :           U       U          100 (D)
Shutdown multiplier :     U       U          1 (D)
Auto-restore multiplier :  U       U          10 (D)
```

```
RP/0/RP0/CPU0:ios(config)#interface hundredGigE 0/2/0/1 priority-flow-control $
RP/0/RP0/CPU0:ios(config)#commit
```

```
RP/0/RP0/CPU0:ios(config)#do show controllers hundredGigE 0/2/0/1 priority-flo$
```

```
Priority flow control information for interface HundredGigE0/2/0/1:
```

```
Priority flow control watchdog configuration:
```

```
(D) : Default value
```

```
U: Unconfigured
```

```
-----
Configuration Item                Global   Interface  Effective
-----
PFC watchdog state :              Enabled  Disabled   Disabled
Poll interval :                   U       U          100 (D)
Shutdown multiplier :             U       U          1 (D)
Auto-restore multiplier :         U       U          10 (D)
-----
```

Related Topics

- [Priority Flow Control Watchdog Overview, on page 79](#)

Global Statistics Counters for Priority Flow Control and Priority Flow Control Watchdog

Table 19: Feature History Table

Feature Name	Release Information	Feature Description
Global Statistics Counters for Priority Flow Control and Priority Flow Control Watchdog	Release 7.5.5	<p>You can now view statistics for Priority Flow Control (PFC) and PFC Watchdog for all interfaces in a consolidated, compact, tabular, and easy-to-read format. We've also made the display of these global statistics faster by ensuring data is collected from all line cards for their interfaces and cumulatively sent to the local statistics infrastructure from where the show commands collect the data.</p> <p>Previously, you could view statistics only per interface for PFC and PFC Watchdog.</p> <p>This feature introduces these changes:</p> <ul style="list-style-type: none"> • CLI: <ul style="list-style-type: none"> • show controllers all priority-flow-control statistics location all • show controllers all priority-flow-control watchdog-stats location all

Previously, you could view statistics for Priority Flow Control (PFC) and PFC Watchdog only per interface by running the following commands:

- For PFC: **show controllers priority-flow-control statistics**
- For PFC Watchdog: **show controllers priority-flow-control watchdog statistics**

With Release 7.5.5, you can view the statistics for PFC and PFC Watchdog for all interfaces by running the following commands.

- For PFC: **show controllers all priority-flow-control statistics location all**
- For PFC Watchdog: **show controllers all priority-flow-control watchdog-stats location all**

Global Statistics Counters for Priority Flow Control and Priority Flow Control Watchdog: Benefits

The statistics display in a consolidated, compact, and easy-to-read format. Plus, unlike single interface statistics, where the show commands get the data from each interface, we've enabled the cumulative data collection from all line cards for their interfaces in the local statistics infrastructure. This displays the statistics faster because the **show controllers all priority-flow-control statistics location all** and **show controllers all priority-flow-control watchdog-stats location all** commands fetch the data from this infrastructure.

Global Statistics Counters for Priority Flow Control and Priority Flow Control Watchdog: Guidelines and Limitations

- You can still view statistics per interface for PFC and PFC Watchdog by running the **show controllers priority-flow-control statistics** and **show controllers priority-flow-control watchdog statistics** commands, respectively.
- Running the **clear controller priority-flow-control statistics** and **clear controllers priority-flow-control watchdog statistics** commands also clears the counters for the corresponding interfaces that the **show controllers all priority-flow-control statistics location all** and **show controllers all priority-flow-control watchdog-stats location all** commands display.

View Global Statistics for Priority Flow Control and Priority Flow Control Watchdog

Run the **show controllers all priority-flow-control statistics location all** command to view the following PFC details for all interfaces:

- **Interface:** Interface name
- **TC:** Traffic Class
- **Rx PFC:** PFC frames received for the traffic class
- **Tx PFC:** PFC frames transmitted for the traffic class

- **Rx Dropped:** Received PFC frames that are dropped. (0 indicates no frames were dropped, while NA indicates the information is 'Not Available'— the line card hasn't sent the data.)

```
Router#show controllers all priority-flow-control statistics location all
```

Interface	TC	Rx PFC	TxPFC	RxDropped
FourHundredGigE0/0/0/0	0	0	0	NA
FourHundredGigE0/0/0/0	7	0	0	NA
FourHundredGigE0/0/0/0	all	0	0	0
FourHundredGigE0/0/0/1	0	0	0	NA
FourHundredGigE0/0/0/1	7	0	0	NA
FourHundredGigE0/0/0/1	all	0	0	0
FourHundredGigE0/0/0/2	0	0	0	NA
FourHundredGigE0/0/0/2	7	0	0	NA
FourHundredGigE0/0/0/2	all	0	0	0
FourHundredGigE0/0/0/3	0	0	0	NA

Run the **show controllers all priority-flow-control watchdog-stats location all** to view the following PFC Watchdog details for all interfaces:

- **Interface:** interface name
- **TC:** traffic class
- **Watchdog events:** the number of times the Watchdog module receives notifications that excess PFC frames are received.
- **Shutdown events:** The number of times the PFC Watchdog moves the queue to the shutdown state.
- **Auto restore Events:** The number of times the PFC Watchdog sets the auto-restore timer.
- **Total Dropped:** Cumulative VOQ drops + output drops in all the watchdog events from the time you ran the clear controllers priority-flow-control watchdog statistics command to clear the statistics counter.
- **Dropped:** Cumulative VOQ drops + output queue drops in the most recent Watchdog shutdown event.

```
Router#show controllers all priority-flow-control watchdog-stats location all
```

Interface	TC	Watchdog	Shutdown	Auto Restore	SAR	SAR
Instantaneous	Total	Events	Dropped	Events	Events	Events
FourHundredGigE0/0/0/0	0	0	0	0	0	0
0	0	0	0	0	0	0
FourHundredGigE0/0/0/0	7	0	0	0	0	0
0	0	0	0	0	0	0
FourHundredGigE0/0/0/1	0	0	0	0	0	0
0	0	0	0	0	0	0
FourHundredGigE0/0/0/1	7	0	0	0	0	0
0	0	0	0	0	0	0
FourHundredGigE0/0/0/2	0	0	0	0	0	0
0	0	0	0	0	0	0
FourHundredGigE0/0/0/2	7	0	0	0	0	0
0	0	0	0	0	0	0
FourHundredGigE0/0/0/3	0	0	0	0	0	0
0	0	0	0	0	0	0
FourHundredGigE0/0/0/3	7	0	0	0	0	0
0	0	0	0	0	0	0

```
FourHundredGigE0/0/0/4 0 0 0 0 0
0 0 0 0 0
```



Note Disregard the **SAR Events** and **SAR Instantaneous Events** entries because those numbers have no bearing on your operations.

PFC Watchdog Statistics to Monitor Traffic Drops

Table 20: Feature History Table

Feature Name	Release Information	Feature Description
PFC Watchdog Statistics to Monitor Traffic Drops	Release 7.5.4	<p>We have enhanced the PFC Watchdog counters to ensure you get an accurate view of the packets dropped across queues. You now get separate statistic counters for the number of packets dropped and the total number of packets dropped.</p> <p>By providing deeper insights into Watchdog packet drops, we aim to help you plan your traffic flows better and avoid potential PFC storms.</p> <p>Earlier releases didn't support counters for packet drops.</p> <p>You can view the counters using:</p> <ul style="list-style-type: none"> • YANG data model (at Github under the 753 folder): Cisco-IOS-XR-ofa-npu-pfc-oper • CLI: show controllers priority-flow-control watchdog-stats

Here's how the PFC Watchdog works when it detects that a queue has stalled:

1. It begins dropping packets that arrive at the stalled queue. (See [Priority Flow Control Watchdog Overview, on page 79](#) for details.)
2. After the congestion clears and no more PFC frames arrive at the queue, the Watchdog restores the queue and flushes all packets that are stuck in the Virtual Output Queue (VOQ) and the output queue.

This functionality with an enhanced statistics view provides two new counters for all the dropped packets.

PFC Watchdog Statistics for Two Counters

With Release 7.5.4, the PFC Watchdog statistics account for two counters:

- Dropped Packets = Cumulative VOQ drops + output queue drops in the most recent Watchdog shutdown event
- Total Dropped Packets = Cumulative VOQ drops + output drops in all the watchdog events from the time you ran the `clear controllers priority-flow-control watchdog statistics` command to clear the statistics counter.

Thus, when you run the `show controllers priority-flow-control watchdog-stats` command, the output displays the statistics for dropped packets and total dropped packets.

PFC Watchdog Statistics Benefits

- With the additional insights the two counters provide, you can study traffic drop patterns across longer timelines.
- You can analyze traffic drops for every congestion recovery that a queue makes.
- With access to such deep analysis, you can plan your traffic flows better. You can adjust your traffic flow rates and prevent potential PFC storms based on previous drop patterns and trends.

View PFC Watchdog Statistics

```
Router#show controllers hundredGigE 0/1/0/43 priority-flow-control watchdog-stats
```

```
Priority flow control information for interface HundredGigE0/1/0/43:
```

```
Priority flow control watchdog statistics:
```

```
SAR: Auto restore and shutdown
```

Traffic Class	:	0	1	2	3	4	5	6
7								
Watchdog Events	:	0	0	0	3	3	0	0
0								
Shutdown Events	:	0	0	0	3	3	0	0
0								
Auto Restore Events	:	0	0	0	3	3	0	0
0								
SAR Events	:	0	0	0	3510	3510	0	0
0								
SAR Instantaneous Events	:	0	0	0	1172	1172	0	0
0								
Total Dropped Packets	:	0	0	0	941505767	941488166	0	0
0								
Dropped Packets	:	0	0	0	314855466	314887161	0	0
0								



Note Disregard the **SAR Events** and **SAR Instantaneous Events** entries because those numbers have no bearing on your operations.

Detect High Bandwidth Memory Congestion

Table 21: Feature History Table

Feature Name	Release Information	Feature Description
Detect High Bandwidth Memory Congestion	Release 7.5.3	<p>We provide detailed insights into congestion on the High Bandwidth Memory (HBM), such as the devices on which congestion has occurred, the time stamps, and when the device returned to its normal state. With such details, you can investigate the cause of congestion and identify the source ports causing congestion for future preventive actions.</p> <p>You must configure PFC in the buffer-extended mode for this option.</p> <p>The feature introduces the following to enable the option to detect HBM congestion:</p> <ul style="list-style-type: none"> • YANG data model (at Github under the 753 folder): Cisco-IOS-XR-um-8000-hw-module-profile-cfg • CLI: hw-module profile npu memory buffer-extended location bandwidth-congestion-detection enable <p>It also introduces the following to view the congestion and memory usage details:</p> <ul style="list-style-type: none"> • YANG data model (at Github under the 753 folder): Cisco-IOS-XR-8000-platforms-npu-memory-oper • CLI: show controllers npu packet-memory

Here's how congestion build-up in queues drives the usage of memory on your router:

1. Under normal queue conditions, packets egressing from an interface are enqueued into the Shared Memory System (SMS) or internal buffer memory.
2. When congestion occurs, the SMS remains the buffer for packets until the congestion exceeds the per-queue usage criteria. This usage criterion is a combination of the threshold buffer space per VOQ and the age of the packet (defined in units of milliseconds).

3. At this stage, the packets are evicted to the High Bandwidth Memory (HBM) or external buffer memory. This provides additional buffer memory to absorb short bursts of congestion.
4. However, when packets from multiple VOQs are evicted to the HBM, the HBM bandwidth begins to experience congestion. The HBM congestion causes data loss for lossless traffic classes.

It's for insights into these extreme events that we've provided the option to detect congestion in the HBM. This information is meant for post-event analysis and reporting. You can enable this functionality by configuring the command **hw-module profile npu memory buffer-extended bandwidth-congestion-detection enable**. By doing so, you can view information such as:

- Devices on which congestion has occurred and the time stamps
- The current buffer memory usage and the highest memory usage watermark reached since the last reading

Benefits of Detecting High Bandwidth Memory Congestion

When you enable the ability to detect HBM congestion, the **show controllers npu packet-memory** command displays:

- Details of congestion along with the affected devices.
- A snapshot of the maximum memory usage with time stamps.

Using this information, you can conduct post-event analysis such as:

- investigate the cause of the congestion.
- identify the ports and flows that are the source of this congestion.

High Bandwidth Memory Congestion Detection: Guidelines and Limitations

- You must configure PFC in the buffer-extended mode for this functionality.
- There's no requirement for a line card reload after you configure **hw-module profile npu memory buffer-extended bandwidth-congestion-detection enable**.
- This functionality is supported on:
 - Cisco Silicon One Q200-based routers and line cards
 - 8201-32FH routers
 - 88-LC0-48TH-MO line cards
 - 88-LC0-36FH line cards
 - 88-LC0-36FH-M line cards

Configure High Bandwidth Memory Congestion Detection

To enable the detection of HBM congestion, configure the **hw-module profile npu memory buffer-extended bandwidth-congestion-detection enable** command.

```

Router#config
Router(config)#hw-module profile npu buffer-extended location 0/6/CPU0
bandwidth-congestion-detection enable
Router(config)#commit
Router(config)#exit

```



Note There's no requirement to reload the line card after the configuration.

Verification

The following table lists the various verification commands you can run depending on your required information. The outputs for these commands are available after the table.

Table 22: Verification Commands

Type of Information	Commands	Details
Buffer usage	<ul style="list-style-type: none"> show controllers npu packet-memory usage instance all location all show controllers npu packet-memory usage verbose instance all location all 	<ul style="list-style-type: none"> These commands provide data for the current use and the highest watermark reached since the last reading for both SMS and HBM. The refresh interval for the information is 30 seconds.
HBM bandwidth congestion	<ul style="list-style-type: none"> show controllers npu packet-memory congestion instance all location all show controllers npu packet-memory congestion detail instance all location all show controllers npu packet-memory congestion verbose instance all location all 	<ul style="list-style-type: none"> These commands provide data for when the HBM congestion occurred or is about to happen. The output maintains a history of the last 120 events regardless of the elapsed time. The refresh interval for new events to be added is 30 seconds.

Run the **show controllers npu packet-memory usage instance all location all** command to view the following details:

- the timestamp at which data is sampled
- the network processor name (**Device**)
- packet memory usage for that timestamp for SMS (**Buff-int Usage** in units of buffers) and HBM (**Buff-ext Usage** in units of 8 KB blocks)
- highest maximum watermark reached for SMS (**Buff-int Max WM**) and HBM (**Buff-ext Max WM**) since the last reading

```
Router#show controllers npu packet-memory usage instance all location all
HW memory Information For Location: 0/6/CPU0
```

```
-----
Timestamp(msec)          | Device | Buff-int | Buff-int | Buff-ext | Buff-ext
                        | Usage  | Max WM   | Usage    | Max WM
-----
Wed 2022-09-21 01:54:11.154 UTC    0       7         8         0         0
Wed 2022-09-21 01:54:12.154 UTC    0       7         8         0         0
Wed 2022-09-21 01:54:24.023 UTC    1      22        22         0         0
Wed 2022-09-21 01:54:34.088 UTC    2      11        12         0         0
Wed 2022-09-21 01:54:35.088 UTC    2      11        12         0         0
Wed 2022-09-21 01:54:36.088 UTC    2      11        12         0         0
Wed 2022-09-21 01:54:37.088 UTC    2      11        12         0         0
Wed 2022-09-21 01:54:38.089 UTC    2      11        12         0         0
Wed 2022-09-21 01:54:39.089 UTC    2      11        12         0         0
Wed 2022-09-21 01:54:40.089 UTC    2      11        12         0         0
```

In addition, the **show controllers npu packet-memory usage verbose instance all location all** command displays timestamp in milliseconds:

```
RP/0/RP0/CPU0:Router#show controllers npu packet-memory usage verbose instance all location all
```

```
HW memory Information For Location: 0/RP0/CPU0
```

```
* Option 'verbose' formatted data is for internal consumption.
```

```
-----
Timestamp(msec) | Device | Buff-int | Buff-int | Buff-ext | Buff-ext
                | Usage  | Max WM   | Usage    | Max WM
-----
1663958881006   0      2455     2676     637      640
1663958882007   0      2461     2703     635      640
1663958883007   0      2364     2690     635      640
1663958884007   0     71603   75325    3183    18336
1663958885008   0      2458     2852     1275    1279
1663958886008   0      2484     2827     1275    1279
```

Run the **show controllers npu packet-memory congestion instance all location all** command to view if HBM congestion has occurred and the timestamp of the congestion state.

- Each row indicates:
 - the timestamp at which data is sampled
 - the network processor name (**Device**)
 - the event type (**Normal** or **Congest**)
- The data displayed is for the last 120 events, and new events get added every 30 seconds. To view the updated data, re-run the command.
- After 120 events, the latest entry replaces the oldest entry. You can't clear the events from the list.

```
Router#show controllers npu packet-memory congestion instance all location all
HW memory Information For Location: 0/6/CPU0
```

```
-----
Timestamp(msec)          |           Buff-ext           | Device
                        | Event Type                    |
-----
Wed 2022-09-21 02:14:41.709 UTC          Congest          1
```

```

Wed 2022-09-21 02:14:41.959 UTC          Congest      1
Wed 2022-09-21 02:14:42.960 UTC          Congest      1
Wed 2022-09-21 02:14:43.960 UTC          Congest      1
Wed 2022-09-21 02:14:45.210 UTC          Congest      1
Wed 2022-09-21 02:14:45.710 UTC          Congest      1
Wed 2022-09-21 02:14:47.711 UTC          Normal       1
    
```

Run the **show controllers npu packet-memory congestion detail instance all location all** command to view the following details:

- the timestamp at which data is sampled
- the event type (**Normal** or **Congest**)
- the network processor name (**Device**) and the slice number (**Slice**) for that device. Every network processor has a fixed number of slices, and each slice, in turn, has a set number of ports.
- single VOQ buffer and aggregated SMS VOQ buffers
- packet memory usage for that timestamp for SMS (**Buff-int Usage** in units of buffers) and HBM (**Buff-ext Usage** in units of 8 KB blocks)
- highest maximum watermark reached for SMS (**Buff-int Max WM**) and HBM (**Buff-ext Max WM**) since the last reading

```
RRP/0/RP0/CPU0:ios#show controllers npu packet-memory congestion detail instance all location all
```

```
Fri Sep 23 18:49:50.640 UTC
HW memory Information For Location: 0/RP0/CPU0
```

* Option 'detail' formatted data is for internal consumption.

Timestamp (msec)	Device	Slice	VOQ	VOQ-buff	Event Type	Buff-int Max WM	Buff-ext Usage	Buff-ext Max WM	Buff-int Usage	Buff-int Max WM	Evicted-buff int-WM
Fri 2022-09-23 18:42:30.349 UTC		5	534	16011	Congest	70410	34405	34405	63969	70410	63969
Fri 2022-09-23 18:42:31.101 UTC		5	534	0	Normal	2440	0	0	900	2440	0
Fri 2022-09-23 18:42:37.354 UTC		5	534	16011	Congest	70573	34408	34408	63984	70573	63984
Fri 2022-09-23 18:42:38.354 UTC		5	534	0	Normal	2455	0	0	915	2455	0
Fri 2022-09-23 18:42:44.606 UTC		5	534	16011	Congest	70081	34532	34532	64002	70081	64002

Run the **show controllers npu packet-memory congestion verbose instance all location all** command to view the following details:

- timestamp in milliseconds when the data is sampled
- the network processor name (**Device**) and the slice number (**Slice**) for that device. Every network processor has a fixed number of slices, and each slice, in turn, has a set number of ports
- event type, where 0 = single VOQ-based congestion and 1 = single VOQ-based congestion backoff (**VOQ-buff int-WM**), 2 = congestion in aggregated SMS buffers for VOQ and 3 = congestion backoff in aggregated SMS buffers for VOQ (**Evicted-buff int-WM**)
- buffer internal for unicast, which is for information only.

- packet memory usage for that timestamp for SMS (**Buff-int Usage** in units of buffers) and HBM (**Buff-ext Usage** in units of 8 KB blocks)
- highest maximum watermark reached for SMS (**Buff-int Max WM**) and HBM (**Buff-ext Max WM**) since the last reading

```
Router#show controllers npu packet-memory congestion verbose instance all location all
HW memory Information For Location: 0/RP0/CPU0
* Option 'verbose' formatted data is for internal consumption.
```

Timestamp(msec)	Event	Device	Slice	VOQ	VOQ-buff	Evicted-buff	Buff-int
Buff-int	Buff-int	Buff-ext	Buff-ext		int-WM	int-WM	UC-WM
Usage	Max WM	Usage	Max WM				
1663958550349	0	0	5	534	16011	63969	65451
70410	70410	34405	34405				
1663958551101	1	0	5	534	0	0	900
2440	2440	0	0				
1663958557354	0	0	5	534	16011	63984	65493
70573	70573	34408	34408				
1663958558354	1	0	5	534	0	0	915
2455	2455	0	0				
1663958564606	0	0	5	534	16011	64002	65520
70081	70081	34532	34532				
1663958565356	1	0	5	534	0	0	915
2417	2417	0	0				

Available Shared Memory System and High Bandwidth Memory

Table 23: Feature History Table

Feature Name	Release Information	Feature Description
Available Shared Memory System and High Bandwidth Memory Buffers	Release 7.5.5	<p>You can now view buffer availability for Shared Memory System (SMS) and High Bandwidth Memory (HBM) with higher accuracy without any lag between the minimum and maximum watermark readings, especially when the packet buffers are used and released rapidly. This is possible because we've enabled the instantaneous display of available or free SMS and HBM.</p> <p>Previously, you could view details only for the highest watermark readings for SMS and HBM.</p> <p>You must configure PFC in the buffer-extended mode for this option, and this functionality is available only for Cisco Silicon One Q200-based routers and line cards.</p> <p>This functionality modifies the following:</p> <ul style="list-style-type: none"> • CLI: show controllers npu packet-memory • YANG Data Model: Cisco-IOS-XR-8000-platforms-npu-memory-oper

From Release 7.5.5 onwards, you can additionally view details of the buffer available for SMS (**Buff-int-free Min WM**) and HBM (**Buff-ext-free Min WM**) when you run the **show controllers npu packet-memory usage instance all location all** command.

These counters are accounted for as follows:

- Available SMS buffer (**Buff-int-free Min WM**) at a given instant = Maximum SMS buffer - highest maximum watermark reached for SMS (**Buff-int Max WM**)
- Available HBM buffer (**Buff-ext-free Min WM**) at a given instant = Maximum HBM buffer - highest maximum watermark reached for HBM (**Buff-ext Max WM**)

This additional detail, available only for Cisco Silicon One Q200-based routers and line cards, provides even better insights to analyze the congestion affecting buffer occupancy, especially when packet buffers are used and released quickly.

```
Router#show controller npu packet-memory usage instance all location all
HW memory Information For Location: 0/6/CPU0
```

```
-----
Timestamp(msec)           | Device | Buff-int | Buff-int | Buff-ext | Buff-ext |
Buff-int-free | Buff-ext-free           | Usage   | Max WM  | Usage   | Max WM  | Min
WM                   | Min WM
-----
Wed 2023-08-30 23:47:40.918 UTC      0      1518      6668      17154      17656
```

293394	982846						
Wed 2023-08-30 23:47:41.918 UTC		0	1227	5631	16010	16427	
293685	983990						
Wed 2023-08-30 23:47:42.919 UTC		0	1398	8295	15041	15734	
293514	984959						
Wed 2023-08-30 23:47:43.919 UTC		0	1765	8892	14744	15678	
293147	985256						
Wed 2023-08-30 23:47:41.011 UTC		1	10380	12419	37532	38165	
284532	962468						
Wed 2023-08-30 23:47:42.011 UTC		1	10463	11977	37315	38326	
284449	962685						
Wed 2023-08-30 23:47:43.013 UTC		1	9145	12604	37714	38242	
285767	962286						
Wed 2023-08-30 23:47:44.013 UTC		1	10996	13272	37429	38051	
283916	962571						

Global Pause Frames for High Bandwidth Memory Congestion

Table 24: Feature History Table

Feature Name	Release Information	Feature Description
Global Pause Frames for High Bandwidth Memory Congestion	Release 7.5.4	<p>We ensure no packet drops on PFC-enabled queues due to High Bandwidth Memory (HBM) congestion. Such prevention of drops is possible because we have enabled the triggering of global pause frames (X-Off) whenever there's HBM congestion.</p> <p>This functionality is disabled by default. You have the following options to enable it:</p> <ul style="list-style-type: none"> • CLI: hw-module profile npu memory buffer-extended bandwidth-congestion-protect enable • YANG Data Model: Cisco-IOS-XR-um-8000-hw-module-profile-cfg (see GitHub, YANG Data Models Navigator) <p>This feature introduces the show hw-module bandwidth-congestion-protect command to view the status of the global X-Off configuration.</p>

When High Bandwidth Memory congestion occurs (see [Detect High Bandwidth Memory Congestion](#), on page 87 for details), global pause frames (X-Off) are triggered for all PFC-enabled queues, regardless of whether those queues are the aggressor queues, hence the name 'global.' The only queues that don't transmit the X-Off trigger are those that don't receive any traffic. Such action ensures no packet drops on lossless queues, allowing you to meet your traffic bandwidth commitments for specific customers and requirements. It also ensures that the X-Off isn't triggered prematurely, which would have affected uncongested queues as well, causing a drop in performance.

To enable this functionality, configure the **hw-module profile npu memory buffer-extended bandwidth-congestion-protect enable** command.

Global Pause Frames for High Bandwidth Memory Congestion: Guidelines and Limitations

- This functionality isn't supported for the buffer-extended mode where the devices are more than 0.5 km apart.
- Configuring the **hw-module profile npu memory buffer-extended bandwidth-congestion-protect enable** command for line cards where you've configured headroom values exceeding 6144000 bytes could result in a commit error or the feature not being enabled.
- You must reload the line card for the **hw-module profile npu memory buffer-extended bandwidth-congestion-protect enable** command to take effect.
- This functionality is supported on:
 - 88-LC0-36FH line cards
 - 88-LC0-36FH-M line cards

Configure Global Pause Frames for High Bandwidth Memory Congestion

To enable triggering of global pause frames (X-off) whenever there's HBM congestion in the buffer-extended mode, configure the **hw-module profile npu memory buffer-extended bandwidth-congestion-protect enable** command.

```
Router#config
Router(config)#hw-module profile npu memory buffer-extended location 0/1/CPU0
bandwidth-congestion-protect enable
Router(config)#commit
```

Verification

Run the **show hw-module bandwidth-congestion-protect** command to view details about global X-off.

```
RP/0/RP1/CPU0:router#show hw-module bandwidth-congestion-protect location 0/1/CPU0
```

Location	Configured	Applied	Action
0/1/CPU0	Yes	No	Reload

The table lists the various possibilities for the command output based on your activity.

If you...	Configured field displays...	Applied field displays...	Action field displays...
Configure the hw-module profile npu memory buffer-extended command	Yes	No	Reload
Use the no form of the hw-module profile npu memory buffer-extended command after configuring it, but before reloading the line card	No	No	N/A

If you...	Configured field displays...	Applied field displays...	Action field displays...
Configure the hw-module profile npu memory buffer-extended command for a supported variant and reload the line card	Yes	Yes, Active Note Yes indicates that the configuration is programmed to the hardware, Active indicates that the global X-off functionality is active on the hardware.	N/A
Use the no form of the hw-module profile npu memory buffer-extended command when it is active, and commit the no form but don't reload the line card	No Note At this stage, the output displays the user action and not the hardware status.	No Note At this stage, the output displays the user action and not the hardware status.	Reload
Reload the line card after committing the no form of the hw-module profile npu memory buffer-extended command	No Note At this stage, the output displays the hardware status.	No Note At this stage, the output displays the hardware status.	N/A



CHAPTER 7

Congestion Management

- [Congestion Management Overview, on page 97](#)
- [Low-Latency Queueing with Strict Priority Queueing, on page 97](#)
- [Committed Bursts, on page 107](#)
- [Excess Bursts, on page 108](#)
- [Two-Rate Policer Details, on page 109](#)
- [References for Modular QoS Management, on page 109](#)

Congestion Management Overview

Congestion management features allow you to control congestion by determining the order in which a traffic flow (or packets) is sent out an interface based on priorities assigned to packets. Congestion management entails the creation of queues, assignment of packets to those queues based on the classification of the packet, and scheduling of the packets in a queue for transmission.

The types of traffic regulation mechanisms supported are:

- [Low-Latency Queueing with Strict Priority Queueing, on page 97](#)
- [Traffic Shaping, on page 100](#)
- [Traffic Policing, on page 102](#)

Low-Latency Queueing with Strict Priority Queueing

The Low-Latency Queueing (LLQ) feature brings strict priority queueing (PQ) to the CBWFQ scheduling mechanism. Priority Queueing (PQ) in strict priority mode ensures that one type of traffic is sent, possibly at the expense of all others. For PQ, a low-priority queue can be detrimentally affected, and, in the worst case, never allowed to send its packets if a limited amount of bandwidth is available or the transmission rate of critical traffic is high. Strict PQ allows delay-sensitive data, such as voice, to be de-queued and sent before packets in other queues are de-queued.

Configure Low Latency Queueing with Strict Priority Queueing

Configuring low latency queueing (LLQ) with strict priority queueing (PQ) allows delay-sensitive data such as voice to be de-queued and sent before the packets in other queues are de-queued.

Guidelines

- Only priority level 1 to 7 is supported, with 1 being the highest priority and 7 being the lowest. However, the default CoSQ 0 has the lowest priority among all.
- Egress policing is not supported. Hence, in the case of strict priority queueing, there are chances that the other queues do not get serviced.
- You can configure **shape average** and **queue-limit** commands along with **priority**.
- You can configure an optional shaper to cap the maximum traffic rate. This is to ensure the lower priority traffic classes are not starved of bandwidth.

Configuration Example

You have to accomplish the following to complete the LLQ with strict priority queueing:

1. Creating or modifying a policy-map that can be attached to one or more interfaces
2. Specifying the traffic class whose policy has to be created or changed
3. Specifying priority to the traffic class
4. (Optional) Shaping the traffic to a specific bit rate
5. Attaching the policy-map to an output interface

```
Router# configure
Router(config)# policy-map test-priority-1
Router(config-pmap)# class qos1
Router(config-pmap-c)# priority level 7
Router(config-pmap-c)# exit
Router(config-pmap)# exit
Router(config)# interface HundredGigE 0/6/0/18
Router(config-if)# service-policy output test-priority-1
Router(config-if)# no shutdown
Router(config-if)# commit
```

Running Configuration

```
policy-map strict-priority
class tc7
  priority level 1
  queue-limit 75 mbytes
!
class tc6
  priority level 2
  queue-limit 75 mbytes
!
class tc5
  priority level 3
  queue-limit 75 mbytes
!
class tc4
  priority level 4
  queue-limit 75 mbytes
!
class tc3
  priority level 5
```

```

    queue-limit 75 mbytes
  !
  class tc2
    priority level 6
    queue-limit 75 mbytes
  !
  class tc1
    priority level 7
    queue-limit 75 mbytes
  !
  class class-default
    queue-limit 75 mbytes
  !
end-policy-map
!

class-map match-any tc1
match traffic-class 1
end-class-map
!
class-map match-any tc2
match traffic-class 2
end-class-map
!
class-map match-any tc3
match traffic-class 3
end-class-map
!
class-map match-any tc4
match traffic-class 4
end-class-map
!
class-map match-any tc5
match traffic-class 5
end-class-map
!
class-map match-any tc6
match traffic-class 6
end-class-map
!
class-map match-any tc7
match traffic-class 7
end-class-map
!

interface HundredGigE0/6/0/18
 service-policy input 100g-s1-1
 service-policy output test-priority-1
!
```

Verification

```
Router# show qos int hundredGigE 0/6/0/18 output
```

```
NOTE:- Configured values are displayed within parentheses
Interface HundredGigE0/6/0/18 ifh 0x3000220 -- output policy
NPU Id:                               3
Total number of classes:                3
Interface Bandwidth:                    100000000 kbps
VOQ Base:                               11176
```

```

VOQ Stats Handle:          0x88550ea0
Accounting Type:          Layer1 (Include Layer 1 encapsulation and above)
-----
Level1 Class (HP7)       = qos-1
Egressq Queue ID        = 11177 (HP7 queue)
TailDrop Threshold      = 125304832 bytes / 10 ms (default)
WRED not configured for this class

Level1 Class (HP6)       = qos-2
Egressq Queue ID        = 11178 (HP6 queue)
TailDrop Threshold      = 125304832 bytes / 10 ms (default)
WRED not configured for this class

Level1 Class             = class-default
Egressq Queue ID        = 11176 (Default LP queue)
Queue Max. BW.          = 101803495 kbps (default)
Queue Min. BW.          = 0 kbps (default)
Inverse Weight / Weight = 1 (BWR not configured)
TailDrop Threshold      = 1253376 bytes / 10 ms (default)
WRED not configured for this class

```

Related Topics

- [Congestion Management Overview, on page 97](#)
- [Traffic Shaping, on page 100](#)

Traffic Shaping

Traffic shaping allows you to control the traffic flow exiting an interface to match its transmission to the speed of the remote target interface and ensure that the traffic conforms to policies contracted for it. Traffic adhering to a particular profile can be shaped to meet downstream requirements, thereby eliminating bottlenecks in topologies with data-rate mismatches.



Note Traffic shaping is supported only in egress direction.

Configure Traffic Shaping

The traffic shaping performed on outgoing interfaces is done at the Layer 1 level and includes the Layer 1 header in the rate calculation.

Guidelines

- Only egress traffic shaping is supported.
- It is mandatory to configure all the eight qos-group classes (including class-default) for the egress policies.
- You can configure **shape average** command along with **priority** command.
- It is recommended that you configure all the eight <traffic-class classes> (including **class-default**) for the egress policies. A limited number of < traffic-class class> combinations are supported, but unicast and multicast traffic-class may not be handled consistently. Hence, such configurations are recommended, when the port egress has only unicast traffic.

Configuration Example

You have to accomplish the following to complete the traffic shaping configuration:

1. Creating or modifying a policy-map that can be attached to one or more interfaces
2. Specifying the traffic class whose policy has to be created or changed
3. Shaping the traffic to a specific bit rate
4. Attaching the policy-map to an output interface

```
Router# configure
Router(config)# policy-map egress_policy1
Router(config-pmap)# class c5
Router(config-pmap-c)# shape average percent 40
Router(config-pmap-c)# exit
Router(config-pmap)# exit
Router(config)# interface HundredGigE 0/1/0/0
Router(config-if)# service-policy output egress_policy1
Router(config-if)# commit
```

Running Configuration

```
policy-map egress_policy1
  class c5
    shape average percent 40
  !
  class class-default
  !
end-policy-map
!

interface HundredGigE0/6/0/18
  service-policy input 100g-s1-1
  service-policy output egress_policy1
!
```

Verification

```
Router# show qos interface hundredGigE 0/6/0/18 output
```

```
NOTE:- Configured values are displayed within parentheses
Interface HundredGigE0/6/0/18 ifh 0x3000220 -- output policy
NPU Id: 3
Total number of classes: 2
Interface Bandwidth: 100000000 kbps
VOQ Base: 11176
VOQ Stats Handle: 0x88550ea0
Accounting Type: Layer1 (Include Layer 1 encapsulation and above)
-----
Level1 Class = c5
Egressq Queue ID = 11177 (LP queue)
Queue Max. BW. = 40329846 kbps (40 %)
Queue Min. BW. = 0 kbps (default)
Inverse Weight / Weight = 1 (BWR not configured)
Guaranteed service rate = 40000000 kbps
TailDrop Threshold = 50069504 bytes / 10 ms (default)
```

```

WRED not configured for this class

Level1 Class                = class-default
Egressq Queue ID           = 11176 (Default LP queue)
Queue Max. BW.             = 101803495 kbps (default)
Queue Min. BW.            = 0 kbps (default)
Inverse Weight / Weight    = 1 (BWR not configured)
Guaranteed service rate    = 50000000 kbps
TailDrop Threshold         = 62652416 bytes / 10 ms (default)
WRED not configured for this class

```

Related Topics

- [Congestion Management Overview, on page 97](#)

Traffic Policing

Traffic policing allows you to control the maximum rate of traffic sent or received on an interface and to partition a network into multiple priority levels or class of service (CoS). Traffic policing manages the maximum rate of traffic through a token bucket algorithm. The token bucket algorithm uses user-configured values to determine the maximum rate of traffic allowed on an interface at a given moment in time. The token bucket algorithm is affected by all traffic entering or leaving the interface (depending on where the traffic policy with traffic policing is configured) and is useful in managing network bandwidth in cases where several large packets are sent in the same traffic stream. By default, the configured bandwidth value takes into account the Layer 2 encapsulation that is applied to traffic leaving the interface.

The router supports the following traffic policing mode:

- Single-Rate Two-Color (SR2C) in color-blind mode. See [Single-Rate Policer, on page 103](#).

Restrictions

- 1R3C policers are not supported.
- Up to eight policers are supported per ingress policy.
- Policers are allocated in multiples of 2, so any request for allocating odd number of policers is internally rounded up by a factor of one.
- Only one conditional marking action is supported. You can set discard class to 0/1, that is used to set either virtual output queueing (VOQ) limits or Random Early Detection (RED) profiles.
- If you configure discard-class explicitly at the class level and not under policer, then the explicit mark action is applied to all the transmitted policer traffic.
- Egress policing is not supported.

Committed Bursts and Excess Bursts

Unlike a traffic shaper, a traffic policer does not buffer excess packets and transmit them later. Instead, the policer executes a “send or do not send” policy without buffering. Policing uses normal or committed burst (bc) values and excess burst values (be) to ensure that the router reaches the configured committed information rate (CIR). Policing decides if a packet conforms or exceeds the CIR based on the burst values you configure.

Burst parameters are based on a generic buffering rule for routers, which recommends that you configure buffering to be equal to the round-trip time bit-rate to accommodate the outstanding TCP windows of all connections in times of congestion. During periods of congestion, proper configuration of the excess burst parameter enables the policer to drop packets less aggressively.

For more details, see [Committed Bursts, on page 107](#) and [Excess Bursts, on page 108](#).

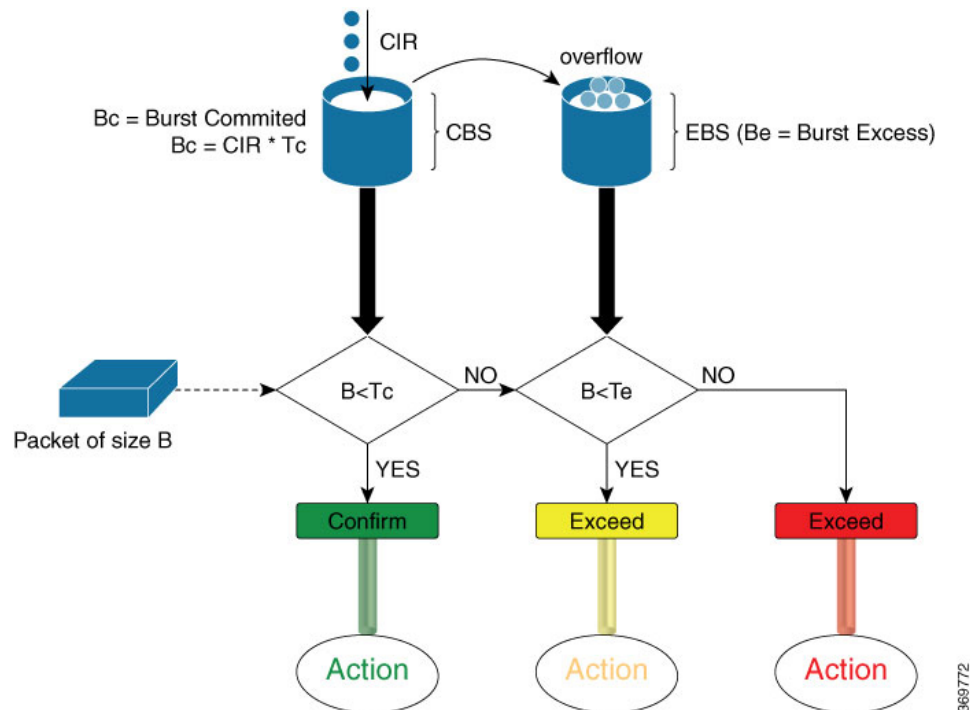
Single-Rate Policer

This section explains the concept of the single-rate two-color policer.

Single-Rate Two-Color Policer

A single-rate two-color (1R2C) policer provides one token bucket with two actions for each packet: a conform action and an exceed action.

Figure 9: Workflow of Single-Rate Two-Color Policing



Based on the committed information rate (CIR) value, the token bucket is updated at every refresh time interval. The T_c token bucket can contain up to the B_c value, which can be a certain number of bytes or a period of time. If a packet of size B is greater than the T_c token bucket, then the packet exceeds the CIR value and a default action is performed. If a packet of size B is less than the T_c token bucket, then the packet conforms and a different default action is performed.

Configure Traffic Policing (Single-Rate Two-Color)

Traffic policing is often configured on interfaces at the edge of a network to limit the rate of traffic entering or leaving the network. The default conform action for single-rate two color policer is to transmit the packet and the default exceed action is to drop the packet. Users cannot modify these default actions.

Configuration Example

You have to accomplish the following to complete the traffic policing configuration:

1. Creating or modifying a policy-map that can be attached to one or more interfaces
2. Specifying the traffic class whose policy has to be created or changed
3. (Optional) Specifying the marking action
4. Specifying the policy rate for the traffic
5. Attaching the policy-map to an input interface

```
Router# configure
Router(config)# policy-map test-police-1R2C
Router(config-pmap)# class dscp1
Router(config-pmap-c)# police rate 10 gbps
Router(config-pmap-c-police)# exit
Router(config-pmap-c)# exit
Router(config-pmap)# exit
Router(config)# interface hundredGigE 0/0/0/18
Router(config-if)# service-policy input test-police-1R2C
Router(config-if)# commit
```

Running Configuration

```
class-map match-any dscp1
  match dscp ipv4 1
end-class-map
!
!
policy-map test-police-1R2C
  class dscp1
    police rate 10 gbps
  !
  !
  class class-default
  !
  !
end-policy-map
!
!
interface HundredGigE0/0/0/8
service-policy input test-police-1R2C
!
```

Verification

```
Router# show qos interface hundredGigE 0/0/0/8 input
NOTE:- Configured values are displayed within parentheses
Interface HundredGigE0/0/0/8 ifh 0xf0001e8 -- input policy
NPU Id: 0
Total number of classes: 2
Interface Bandwidth: 100000000 kbps
Policy Name: test-police-1R2C
Accounting Type: Layer1 (Include Layer 1 encapsulation and above)
-----
Level1 Class = dscp1
```

```

Policer committed rate           = 10000000 kbps (10 gbits/sec)
Policer conform burst            = 1024000 bytes (default)
Policer conform action           = Just TX
Policer exceed action            = DROP PKT

Levell1 Class                    = class-default
Policer not configured for this class

```

Related Topics

- [Traffic Policing, on page 102](#)

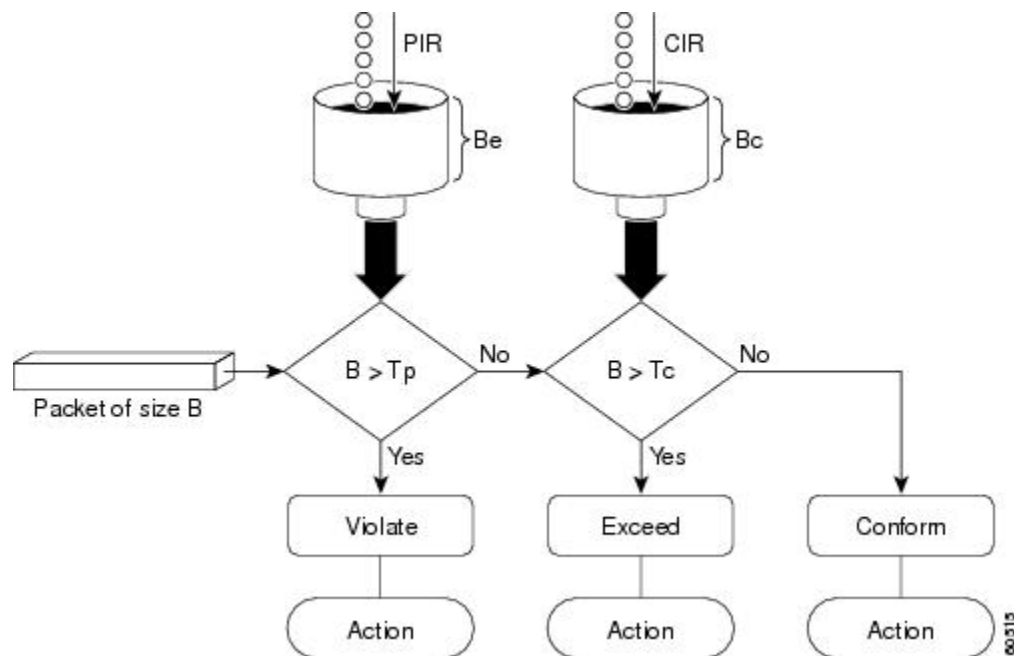
Two-Rate Policer

The two-rate policer manages the maximum rate of traffic by using two token buckets: the committed token bucket and the peak token bucket. The dual-token bucket algorithm uses user-configured values to determine the maximum rate of traffic allowed on a queue at a given moment. In this way, the two-rate policer can meter traffic at two independent rates: the committed information rate (CIR) and the peak information rate (PIR).

The dual-token bucket algorithm provides users with three actions for each packet—a conform action, an exceed action, and an optional violate action. Traffic entering a queue with the two-rate policer configured is placed into one of these categories. The actions are pre-determined for each category. The default conform and exceed actions are to transmit the packet, and the default violate action is to drop the packet.

This figure shows how the two-rate policer marks a packet and assigns a corresponding action to the packet.

Figure 10: Marking Packets and Assigning Actions—Two-Rate Policer



Also, see [Two-Rate Policer Details, on page 109](#).

The router supports Two-Rate Three-Color (2R3C) policer.

Configure Traffic Policing (Two-Rate Three-Color)

The default conform and exceed actions for two-rate three-color (2R3C) policer are to transmit the packet and the default violate action is to drop the packet. Users cannot modify these default actions.

Configuration Example

You have to accomplish the following to complete the two-rate three-color traffic policing configuration:

1. Creating or modifying a policy-map that can be attached to one or more interfaces
2. Specifying the traffic class whose policy has to be created or changed
3. Specifying the packet marking
4. Configuring two rate traffic policing
5. Attaching the policy-map to an input interface

```
Router# configure
Router(config)# policy-map test-police-2R3C
Router(config-pmap)# class dscp1
Router(config-pmap-c)# police rate 10 gbps peak-rate 20 gbps
Router(config-pmap-c-police)# exit
Router(config-pmap-c)# exit
Router(config-pmap)# exit
Router(config)# interface hundredGigE 0/0/0/8
Router(config-if)# service-policy input test-police-2R3C
Router(config-if)# commit
```

Running Configuration

```
class-map match-any dscp1
match dscp ipv4 1
end-class-map
!
!
policy-map test-police-2R3C
class dscp1
police rate 10 gbps peak-rate 20 gbps
!
!
class class-default
!
!
end-policy-map
!
!
interface HundredGigE0/0/0/8
service-policy input test-police-2R3C
!
```

Verification

```
Router# show qos interface hundredGigE 0/0/0/8 input
```

NOTE:- Configured values are displayed within parentheses

```

Interface HundredGigE0/0/0/8 ifh 0xf0001e8 -- input policy
NPU Id:                                0
Total number of classes:                2
Interface Bandwidth:                    100000000 kbps
Policy Name:                            test-police-2R3C
Accounting Type:                        Layer1 (Include Layer 1 encapsulation and above)
-----
Levell Class                            = dscp1
Policer committed rate                   = 10000000 kbps (10 gbits/sec)
Policer peak rate                       = 20000000 kbps (20 gbits/sec)
Policer conform burst                    = 1024000 bytes (default)
Policer exceed burst                     = 2048000 bytes (default)
Policer conform action                   = Just TX
Policer exceed action                    = DROP PKT
Policer violate action                   = DROP PKT

Levell Class                            = class-default
Policer not configured for this class

Router# policy-map interface hundredGigE 0/0/0/8 input

HundredGigE0/0/0/8 input: test-police-2R3C

Class dscp1
  Classification statistics              (packets/bytes)      (rate - kbps)
  Matched                               : 289228439/289228439000    27734775
  Transmitted                            : 56422213/56422213000     5410359
  Total Dropped                          : 232806226/232806226000   22324416
  Policing statistics                    (packets/bytes)      (rate - kbps)
  Policed(conform)                       : 56422213/56422213000    5410359
  Policed(exceed)                        : 56422215/56422215000    5410358
  Policed(violate)                       : 176384011/176384011000  16914058
  Policed and dropped                    : 232806226/232806226000
Class class-default
  Classification statistics              (packets/bytes)      (rate - kbps)
  Matched                               : 61136620/61136620000     0
  Transmitted                            : 61136620/61136620000     0
  Total Dropped                          : 0/0                      0
Policy Bag Stats time: 1570155764000 [Local Time: 10/04/19 02:22:44.000]

```

Related Topics

- [Two-Rate Policer, on page 105](#)

Committed Bursts

The committed burst (bc) parameter of the police command implements the first, conforming (green) token bucket that the router uses to meter traffic. The bc parameter sets the size of this token bucket. Initially, the token bucket is full and the token count is equal to the committed burst size (CBS). Thereafter, the meter updates the token counts the number of times per second indicated by the committed information rate (CIR).

The following describes how the meter uses the conforming token bucket to send packets:

- If sufficient tokens are in the conforming token bucket when a packet arrives, the meter marks the packet green and decrements the conforming token count by the number of bytes of the packet.

- If there are insufficient tokens available in the conforming token bucket, the meter allows the traffic flow to borrow the tokens needed to send the packet. The meter checks the exceeding token bucket for the number of bytes of the packet. If the exceeding token bucket has a sufficient number of tokens available, the meter marks the packet.

Green and decrements the conforming token count down to the minimum value of 0.

Yellow, borrows the remaining tokens needed from the exceeding token bucket, and decrements the exceeding token count by the number of tokens borrowed down to the minimum value of 0.

- If an insufficient number of tokens is available, the meter marks the packet red and does not decrement either of the conforming or exceeding token counts.



Note When the meter marks a packet with a specific color, there must be a sufficient number of tokens of that color to accommodate the entire packet. Therefore, the volume of green packets is never smaller than the committed information rate (CIR) and committed burst size (CBS). Tokens of a given color are always used on packets of that color.

Excess Bursts

The excess burst (be) parameter of the police command implements the second, exceeding (yellow) token bucket that the router uses to meter traffic. The exceeding token bucket is initially full and the token count is equal to the excess burst size (EBS). Thereafter, the meter updates the token counts the number of times per second indicated by the committed information rate (CIR).

The following describes how the meter uses the exceeding token bucket to send packets:

- When the first token bucket (the conforming bucket) meets the committed burst size (CBS), the meter allows the traffic flow to borrow the tokens needed from the exceeding token bucket. The meter marks the packet yellow and then decrements the exceeding token bucket by the number of bytes of the packet.
- If the exceeding token bucket does not have the required tokens to borrow, the meter marks the packet red and does not decrement the conforming or the exceeding token bucket. Instead, the meter performs the exceed-action configured in the police command (for example, the policer drops the packets).

Important Points

- User configurable burst values— committed burst size (CBS) and excess burst size (EBS)—are not supported. Instead, they are derived using user configured rates—committed information rates (CIR) and peak information rates (PIR).
- The router supports two burst values: low (10 kilobytes) and high (1 megabyte). For a lower range of configured values (1.6 MBps to less than 1 GBps), the burst value is 10 KB. For a higher range of configured values (1 GBps to 400 GBps), the burst value is 1 MB.

Two-Rate Policer Details

The committed token bucket can hold bytes up to the size of the committed burst (bc) before overflowing. This token bucket holds the tokens that determine whether a packet conforms to or exceeds the CIR as the following describes:

- A traffic stream is conforming when the average number of bytes over time does not cause the committed token bucket to overflow. When this occurs, the token bucket algorithm marks the traffic stream green.
- A traffic stream is exceeding when it causes the committed token bucket to overflow into the peak token bucket. When this occurs, the token bucket algorithm marks the traffic stream yellow. The peak token bucket is filled as long as the traffic exceeds the police rate.

The peak token bucket can hold bytes up to the size of the peak burst (be) before overflowing. This token bucket holds the tokens that determine whether a packet violates the PIR. A traffic stream is violating when it causes the peak token bucket to overflow. When this occurs, the token bucket algorithm marks the traffic stream red.

For example, if a data stream with a rate of 250 kbps arrives at the two-rate policer, and the CIR is 100 kbps and the PIR is 200 kbps, the policer marks the packet in the following way:

- 100 kbps conforms to the rate
- 100 kbps exceeds the rate
- 50 kbps violates the rate

The router updates the tokens for both the committed and peak token buckets in the following way:

- The router updates the committed token bucket at the CIR value each time a packet arrives at the interface. The committed token bucket can contain up to the committed burst (bc) value.
- The router updates the peak token bucket at the PIR value each time a packet arrives at the interface. The peak token bucket can contain up to the peak burst (be) value.
- When an arriving packet conforms to the CIR, the router takes the conform action on the packet and decrements both the committed and peak token buckets by the number of bytes of the packet.
- When an arriving packet exceeds the CIR, the router takes the exceed action on the packet, decrements the committed token bucket by the number of bytes of the packet, and decrements the peak token bucket by the number of overflow bytes of the packet.
- When an arriving packet exceeds the PIR, the router takes the violate action on the packet, but does not decrement the peak token bucket.

See [Two-Rate Policer](#), on page 105.

References for Modular QoS Management

Read this section for more information on committed bursts, excess bursts, and two-rate policer.



CHAPTER 8

Configure Modular QoS on Link Bundles

- [QoS on Link Bundles, on page 111](#)

QoS on Link Bundles

A bundle is a group of one or more ports that are aggregated together and treated as a single link. The router supports Ethernet interfaces and VLAN interfaces (bundle sub-interfaces) bundles. All QoS features currently supported on physical interfaces, are also supported on all link bundle interfaces. Applying QoS on bundle members is not supported.

Restrictions for Link Bundles

- Only Ethernet link bundling is supported.
- Only relative (percentage-based) shaper values can be configured. You cannot configure absolute values
- A bundle interface can only contain physical interface.
- All links within a single bundle must be configured either to run 802.3ad (LACP) or EtherChannel (non-LACP). Mixed links within a single bundle are not supported.
- MAC accounting is not supported on Ethernet link bundles.
- Maximum number of links supported in each link bundle is 64.
- The maximum number of link bundles supported is 128.

Load Balancing

Load balancing function is a forwarding mechanism to distribute traffic over multiple links based on Layer 3 routing information in the router. Per-destination load balancing is only supported on the router, where the router is allowed to distribute packets over one of the links in the bundle. When the per-destination load balancing is enabled, all packets for a certain source-destination pair goes through the same link, though there are multiple links available. In other words, per-destination load balancing can ensure that packets for a certain source-destination pair could arrive in order.

Layer 3 Load Balancing on Link Bundles

Layer 3 load balancing for link bundles is done on Ethernet Flow Points (EFPs) and is based on the IPv4 source and destination addresses in the packet. When Layer 3 service-specific load balancing is configured, all egress bundles are load balanced based on the IPv4 source and destination addresses. When packets do not have IPv4 addresses, default load-balancing (based on the MAC SA/DA fields in the packet header) is used.

Configure QoS on Link Bundles

QoS is configured on link bundles in the same way that it's configured on individual interfaces.

Guidelines

- When you apply a QoS policy on a bundle (ingress or egress direction), the queuing policy is applied at each member interface. The reference bandwidth that is used to calculate shaper or bandwidth values is applied as per the physical member interface bandwidth.
- If a QoS policy is not applied to a bundle interface, both the ingress and egress traffic use the default queue of the per link member port.
- The shape rate specified in the bundle policy-map is not an aggregate for all bundle members. The shape rate applied to the bundle depends on the load balancing of the links. For example, if a policy map with a shape rate of 10 Mbps is applied to a bundle with two member links, and if the traffic is always load-balanced to the same member link, then an overall rate of 10 Mbps applies to the bundle. However, if the traffic is load-balanced evenly between the two links, the overall shape rate for the bundle becomes 20 Mbps.
- If a member is deleted from a bundle, the total bundle statistics changes because the statistics that belongs to the detached link is lost.
- Add bundle members before you apply the QoS policy to the bundle. This action ensures that the QoS policy applies across the bundle and on the new member. If a QoS policy exists and you add a new member later, the policy may not apply on the bundle if the member interface has an incompatible configuration. This scenario may result in errors.
- **Queue-limit error:**
 - If you experience a queue-limit error, remove the existing QoS policy from the link bundle interfaces. However, this removal may take up to ten minutes to come into effect, depending on the time lapsed since the error displayed. To see the time lapsed, check the **retries performed** information in the subsequent syslog.
 - When you encounter an invalid queue-limit error, the syslog lists only one interface, though the error may occur across multiple interfaces. To know the number of affected interfaces, see **operations** in the syslog.

For example:

```
RP/0/RP0/CPU0:Aug 17 09:56:37.417 PDT: BM-DISTRIB[1156]:
%L2-BM-4-ERR_OP_RETRY_THRESHOLD :
Exceeded threshold for the number of retries for 4 operations on members.
First (and error) FortyGigE0/0/0/34 ('DPA_QOSEA' detected the 'warning'
condition 'Invalid Queue Limit.
Value should be in range 614400-390070272 bytes or equivalent value in usec or
msec based on the service rate of the queue'),
```

```

retries performed for 10 hrs 16 mins

RP/0/RP0/CPU0:Aug 17 10:06:40.388 PDT: BM-DISTRIB[1156]:
%L2-BM-4-ERR_OP_RETRY_THRESHOLD :
Exceeded threshold for the number of retries for 4 operations on members.
First (and error) FortyGigE0/0/0/34 ('DPA_QOSEA' detected the 'warning'
condition 'Invalid Queue Limit.
Value should be in range 614400-390070272 bytes or equivalent value in
usec or msec based on the service rate of the queue'),
retries performed for 10 hrs 26 mins

```

- The QoS policy applied on bundle is inherited to all its member links and the reference bandwidth used to calculate shaper/bandwidth is applied as per the physical member interface bandwidth, and not the bundle as a whole.
- Classification resources (TCAM/MAP) and counters are allocated per bundle interface, not per member.
- The policing policy can be applied on bundle main interface or sub interface, but not on both.
- Policer committed information rates (CIR) and peak information rates (PIR) can be configured only in percent units.
- Policers and policer counters are allocated per bundle interface but not per member.
- The queuing policy can be applied on bundle main or bundle sub interface, but not on both.
- Queuing resources (VoQs) and counters are allocated per bundle member interface.
- The classification and marking policies can be applied on the bundle main or the bundle sub interface, but not on both.

Configuration Example

You have to accomplish the following to complete the QoS configuration on link bundles:

1. Creating a class-map
2. Creating a policy-map and specifying the respective class-map
3. Specifying the action type for the traffic
Refer [Attach a Traffic Policy to an Interface, on page 17](#) for details on step 1, 2 and 3.
4. Creating a link bundle
5. Applying traffic policy to the link bundle

```

/* Configure Ether-Bundle and apply traffic policy */
Router(config)# interface Bundle-Ether 12000
Router(config-if)# mtu 9100
Router(config-if)# service-policy input ingress
Router(config-if)# service-policy output egress
Router(config-if)# ipv4 address 100.12.0.0 255.255.255.254
Router(config-if)# bundle maximum-active links 64
Router(config-if)# commit

```

Running Configuration

This example shows how a queuing policy is applied on an Ethernet link bundle. The policy is applied to all interfaces that are members of the Ethernet link bundle.

```
/* Policy-map */

policy-map ingress
  class inet4-classifier-af1
    set qos-group 1
  !
  class inet4-classifier-af2
    set qos-group 2
  !
  class inet4-classifier-af3
    set qos-group 3
  !
  class inet4-classifier-af4
    set qos-group 4
  !
  class inet4-classifier-be1
    set qos-group 5
  !
  class inet4-classifier-nc1
    set qos-group 6
  !
  class class-default
  !
end-policy-map
!

/* Ether Bundle */
interface Bundle-Ether12000
  mtu 9100
  service-policy input ingress
  service-policy output egress
  ipv4 address 100.12.0.0 255.255.255.254
  load-interval 30

!
```

Verification

- Verify that the bundle status is UP.

```
router# show bundle bundle-ether 1200
Wed Dec 16 19:55:49.974 PST

Bundle-Ether12000
  Status: Up
  Local links <active/standby/configured>: 35 / 0 / 35
  Local bandwidth <effective/available>: 3500000000 (3500000000) kbps
  MAC address (source): ea3b.745f.c4b0 (Chassis pool)
  Inter-chassis link: No
  Minimum active links / bandwidth: 1 / 1 kbps
  Maximum active links: 64
  Wait while timer: 2000 ms
  Load balancing: Default
  LACP: Operational
  Flap suppression timer: Off
  Cisco extensions: Disabled
```

```

Non-revertive:           Disabled
mLACP:                  Not configured
IPv4 BFD:               Not configured

```

Port	Device	State	Port ID	B/W, kbps
Hu0/4/0/0	Local	Active	0x8000, 0x0009	100000000
Link is Active				
Hu0/4/0/1	Local	Active	0x8000, 0x000a	100000000
Link is Active				
- - -				
- - -				
Hu0/4/0/35	Local	Active	0x8000, 0x002b	100000000
Link is Active				

- Verify the bundle statistics:

```

router# show policy-map interface bundle-ether 12000

Bundle-Ether12000 input: ingress

Class inet4-classifier-af1
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          : 4647401962/21236124455654    26403040
  Transmitted                       : 4647401962/21236124455654    26403040
  Total Dropped                     : 0/0                      0
Class inet4-classifier-af2
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          : 4502980177/20576584333939    25571493
  Transmitted                       : 4502980177/20576584333939    25571493
  Total Dropped                     : 0/0                      0
Class inet4-classifier-af3
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          : 4647404125/21236213667880    26389086
  Transmitted                       : 4647404125/21236213667880    26389086
  Total Dropped                     : 0/0                      0
Class inet4-classifier-af4
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          : 9291188840/42456120548683    52771168
  Transmitted                       : 9291188840/42456120548683    52771168
  Total Dropped                     : 0/0                      0
Class inet4-classifier-bel
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          : 4647413429/21235847852686    26393414
  Transmitted                       : 4647413429/21235847852686    26393414
  Total Dropped                     : 0/0                      0
Class inet4-classifier-nc1
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          : 9294887621/42473100149807    52778258
  Transmitted                       : 9294887621/42473100149807    52778258
  Total Dropped                     : 0/0                      0

Class class-default
  Classification statistics          (packets/bytes)    (rate - kbps)
  Matched                          : 0/0                      0
  Transmitted                       : 0/0                      0
  Total Dropped                     : 0/0                      0

Bundle-Ether12000 output: egress

Class cl
  Classification statistics          (packets/bytes)    (rate - kbps)

```

```

    Matched          :          16665494532/75878118942463      8760591
    Transmitted      :          16655834643/75834136022017      8760591
    Total Dropped    :                9659889/43982920446        0
  Queueing statistics
    Queue ID         : None (Bundle)
    Taildropped(packets/bytes) : 9659889/43982920446
Class c2
  Classification statistics          (packets/bytes)      (rate - kbps)
    Matched          :          16665421959/75877849543188      8718687
    Transmitted      :          16665421959/75877849543188      8718687
    Total Dropped    :                0/0                          0
  Queueing statistics
    Queue ID         : None (Bundle)
    Taildropped(packets/bytes) : 0/0
Class c3
  Classification statistics          (packets/bytes)      (rate - kbps)
    Matched          :          16665247833/75877509455458      8703470
    Transmitted      :          16665187414/75877234624197      8703470
    Total Dropped    :                60419/274831261           0
  Queueing statistics
    Queue ID         : None (Bundle)
    Taildropped(packets/bytes) : 60419/274831261
Class c4
  Classification statistics          (packets/bytes)      (rate - kbps)
    Matched          :          33330896131/151755393012945      17470745
    Transmitted      :          33330745421/151754709368565      17470745
    Total Dropped    :                150710/683644380           0
  Queueing statistics
    Queue ID         : None (Bundle)
    Taildropped(packets/bytes) : 150710/683644380
Class c5
  Classification statistics          (packets/bytes)      (rate - kbps)
    Matched          :          16878910340/76849791869834      8833394
    Transmitted      :          16878849464/76849514633309      8833394
    Total Dropped    :                60876/277236525            0
  Queueing statistics
    Queue ID         : None (Bundle)
    Taildropped(packets/bytes) : 60876/277236525
Class c6
  Classification statistics          (packets/bytes)      (rate - kbps)
    Matched          :          33330898844/151756094112925      17456785
    Transmitted      :          33330752668/151755427708382      17456785
    Total Dropped    :                146176/666404543           0
  Queueing statistics
    Queue ID         : None (Bundle)
    Taildropped(packets/bytes) : 146176/666404543
Class c7
  Classification statistics          (packets/bytes)      (rate - kbps)
    Matched          :                244106/79922040              74
    Transmitted      :                244106/79922040              74
    Total Dropped    :                0/0                          0
  Queueing statistics
    Queue ID         : None (Bundle)
    Taildropped(packets/bytes) : 0/0
Class class-default
  Classification statistics          (packets/bytes)      (rate - kbps)
    Matched          :          267075066180/1215993441123215      139917482
    Transmitted      :          267075066180/1215993441123215      139917482
    Total Dropped    :                0/0                          0
  Queueing statistics
    Queue ID         : None (Bundle)
    Taildropped(packets/bytes) : 0/0

```

Related Topics

- [QoS on Link Bundles, on page 111](#)

