



Nokia Service Router Linux

QUALITY OF SERVICE GUIDE RELEASE 21.11

3HE 17916 AAAA TQZZA

Issue 1

December 2021

© 2021 Nokia.

Use subject to Terms available at: www.nokia.com/terms/.

Nokia is committed to diversity and inclusion. We are continuously reviewing our customer documentation and consulting with standards bodies to ensure that terminology is inclusive and aligned with the industry. Our future customer documentation will be updated accordingly.

This document includes Nokia proprietary and confidential information, which may not be distributed or disclosed to any third parties without the prior written consent of Nokia.

This document is intended for use by Nokia's customers ("You"/"Your") in connection with a product purchased or licensed from any company within Nokia Group of Companies. Use this document as agreed. You agree to notify Nokia of any errors you may find in this document; however, should you elect to use this document for any purpose(s) for which it is not intended, You understand and warrant that any determinations You may make or actions You may take will be based upon Your independent judgment and analysis of the content of this document.

Nokia reserves the right to make changes to this document without notice. At all times, the controlling version is the one available on Nokia's site.

No part of this document may be modified.

NO WARRANTY OF ANY KIND, EITHER EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO ANY WARRANTY OF AVAILABILITY, ACCURACY, RELIABILITY, TITLE, NON-INFRINGEMENT, MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE, IS MADE IN RELATION TO THE CONTENT OF THIS DOCUMENT. IN NO EVENT WILL NOKIA BE LIABLE FOR ANY DAMAGES, INCLUDING BUT NOT LIMITED TO SPECIAL, DIRECT, INDIRECT, INCIDENTAL OR CONSEQUENTIAL OR ANY LOSSES, SUCH AS BUT NOT LIMITED TO LOSS OF PROFIT, REVENUE, BUSINESS INTERRUPTION, BUSINESS OPPORTUNITY OR DATA THAT MAY ARISE FROM THE USE OF THIS DOCUMENT OR THE INFORMATION IN IT, EVEN IN THE CASE OF ERRORS IN OR OMISSIONS FROM THIS DOCUMENT OR ITS CONTENT.

Copyright and trademark: Nokia is a registered trademark of Nokia Corporation. Other product names mentioned in this document may be trademarks of their respective owners.

The registered trademark Linux® is used pursuant to a sublicense from the Linux Foundation, the exclusive licensee of Linus Torvalds, owner of the mark on a worldwide basis.

© 2021 Nokia.

Table of contents

1 About this guide.....	5
1.1 What's new.....	5
1.2 Precautionary and information messages.....	5
1.3 Conventions.....	6
2 Quality of service overview.....	7
2.1 How QoS works for transit traffic.....	7
2.2 How QoS works for VXLAN traffic.....	9
2.3 How QoS works for router-terminated traffic.....	10
2.4 How QoS works for router-originated traffic.....	10
3 QoS configuration.....	13
3.1 DSCP classifier policies configuration for input traffic.....	13
3.1.1 Configuring DSCP classifier policies.....	13
3.1.2 Using a DSCP classifier for VXLAN traffic.....	14
3.2 DSCP rewrite-rule policies configuration for output traffic.....	14
3.2.1 Configuring DSCP rewrite-rule policies.....	14
3.3 Queue templates configuration.....	15
3.3.1 Configuring queue templates.....	15
3.3.2 Queue depth (maximum burst size).....	16
3.3.2.1 Configuring queue depth (maximum burst size).....	16
3.3.3 WRED slope.....	16
3.3.3.1 Configuring a WRED slope (7250 IXR).....	17
3.3.3.2 Configuring a WRED slope (7220 IXR).....	17
3.3.4 ECN slope.....	18
3.3.4.1 Configuring an ECN slope (7250 IXR).....	18
3.3.4.2 Configuring an ECN slope (7220 IXR).....	18
3.3.5 Queue utilization thresholds.....	19
3.3.5.1 Configuring queue utilization thresholds on 7250 IXR systems.....	19
3.3.5.2 Configuring queue utilization thresholds on 7220 IXR-D2 and D3 systems.....	20
3.3.5.3 Configuring queue utilization thresholds on 7220 IXR-H2 and H3 systems.....	21
3.4 DSCP classifier policy application to subinterfaces.....	22
3.4.1 Applying a DSCP classifier policy to input traffic (7250 IXR).....	22
3.4.2 Applying a DSCP classifier policy to input traffic (7220 IXR).....	23

3.5 Rewrite-rule policy application to subinterfaces.....	23
3.5.1 Applying a rewrite-rule policy to output traffic (7250 IXR).....	23
3.5.2 Applying a rewrite-rule policy to output traffic (7220 IXR).....	24
3.6 Output queue scheduling.....	24
3.6.1 Configuring strict priority (7250 IXR).....	25
3.6.2 Configuring strict priority (7220 IXR-D2 and D3 or 7220 IXR-H2 and H3).....	25
3.6.3 Configuring WRR (7250 IXR).....	26
3.6.4 Configuring WRR (7220 IXR-D2 and D3 or 7220 IXR-H2 and H3).....	26
3.6.5 Configuring forwarding class peak rate.....	27
4 MPLS QoS overview.....	28
4.1 Ingress LER.....	28
4.2 Transit LSR.....	28
4.3 PHP LSR.....	29
4.4 Egress LER.....	29
4.5 Default MPLS traffic-class classifier policy.....	30
5 MPLS QoS configuration.....	31
5.1 Configuring MPLS traffic-class policy.....	31
5.2 Applying MPLS traffic-class policy to input traffic.....	31
5.3 Configuring MPLS rewrite rules.....	31
5.4 Applying MPLS rewrite rules to output traffic.....	32
6 Buffer utilization display.....	33
6.1 Displaying buffer utilization.....	33
7 Displaying QoS statistics.....	35
7.1 Clearing QoS statistics.....	43
7.2 QoS profile resource usage.....	44
7.2.1 Displaying QoS profile resource usage on a 7250 IXR system.....	44

1 About this guide

This document describes configuration details for the Quality of Service (QoS) feature set used with the Nokia Service Router Linux (SR Linux).

This document is intended for network technicians, administrators, operators, service providers, and others who need to understand how the router is configured.



Note:

This manual covers the current release and may also contain some content that will be released in later maintenance loads. See the *SR Linux Release Notes* for information on features supported in each load.

1.1 What's new

Topic	Location
This is a new document. Previous QoS information was found in the <i>SR Linux Configuration Basics</i> and has now moved to this new document.	
MPLS QoS	MPLS QoS overview MPLS QoS configuration

1.2 Precautionary and information messages

The following are information symbols used in the documentation.



DANGER: Danger warns that the described activity or situation may result in serious personal injury or death. An electric shock hazard could exist. Before you begin work on this equipment, be aware of hazards involving electrical circuitry, be familiar with networking environments, and implement accident prevention procedures.



WARNING: Warning indicates that the described activity or situation may, or will, cause equipment damage, serious performance problems, or loss of data.



Caution: Caution indicates that the described activity or situation may reduce your component or system performance.



Note: Note provides additional operational information.



Tip: Tip provides suggestions for use or best practices.

1.3 Conventions

Nokia SR Linux documentation uses the following command conventions.

- **Bold** type indicates a command that the user must enter.
- Input and output examples are displayed in `Courier` text.
- An open right-angle bracket indicates a progression of menu choices or simple command sequence (often selected from a user interface). Example: **start > connect to**.
- Angle brackets (< >) indicate an item that is not used verbatim. For example, for the command **show ethernet <name>**, *name* should be replaced with the name of the interface.
- A vertical bar (|) indicates a mutually exclusive argument.
- Square brackets ([]) indicate optional elements.
- Braces ({ }) indicate a required choice. When braces are contained within square brackets, they indicate a required choice within an optional element.
- *Italic* type indicates a variable.

Generic IP addresses are used in examples. Replace these with the appropriate IP addresses used in the system.

2 Quality of service overview

Quality of Service (QoS) provides an appropriate level of service for packets as they flow inside the switch and between switches in the network. The required level of service depends on the application that generates the flow of packets, and can be defined by the application's sensitivity to packet loss, delay, and jitter.

QoS functionality is supported on the 7250 IXR, 7220 IXR-D2 and D3, and the 7220 IXR-H2 and H3.

You can group packets that require a similar treatment (per-hop behavior) into a Forwarding Class (FC), also known as a behavior aggregate. You can specify up to eight FCs. Traffic is scheduled and can optionally be marked based on its FC.

A configurable drop probability expresses the packet loss sensitivity. Assign a low drop probability to packets that are sensitive to loss. To provide the required congestion management and intelligent discard decisions when congestion occurs, balance the traffic classifications between low, medium, and high drop probability.

2.1 How QoS works for transit traffic

This section describes how QoS applies to transit packets on the SR Linux.

1. Packets are received on a subinterface.
2. Each received packet is classified as belonging to one of eight FCs (fc0 to fc7) and one of three drop probabilities (low, medium, or high).
 - If the configuration of the ingress subinterface refers to a DSCP classifier policy, the policy determines the FC and drop probability level.



Note: If no entry in this policy matches the received DSCP, the assigned FC is fc0 and the drop probability is low. These classifications correspond to a best effort treatment.

- If there is no DSCP classifier policy bound to the ingress subinterface, the default DSCP classifier policy defines the FC and drop probability. See [Table 1: System default DSCP classifier policy](#).



Note: On all VLAN-based subinterfaces, the 802.1p bits are currently ignored for purposes of FC and drop probability classification.

3. A forwarding lookup on the packet determines its egress port.
4. On the 7250 IXR, each unicast packet is associated with a Virtual Output Queue (VOQ) based on the ingress port, egress port, and FC.

On a 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3, the packet is associated directly with an Egress Queue (EGQ) of the egress port, based on the packet's FC and its type (either unicast or multicast).
5. While it waits for its VOQ or EGQ to be serviced, the packet is stored in buffer memory. The total amount of buffer memory varies by platform.

6. The packet is dropped if the buffer memory is close to full or if the Maximum Burst Size (MBS) of the VOQ or EGQ is exceeded.

The MBS is one of the parameters that is configurable in a queue template. When a queue template is applied to a set of queues, all of those queues have the MBS value specified in the template. If the MBS is not specified in a queue template, the default value is platform dependent. The MBS is not a guaranteed allocation of buffer memory.

7. When the packet is Explicit Congestion Notification (ECN)-capable, and ECN is enabled globally with the **qos explicit-congestion-notification** command, and the VOQ or EGQ has an active ECN slope that applies to the packet, then the ECN field may be remarked depending on the current (weighted) queue depth.
- If the current queue depth is below the configured **min-threshold-percent** of the ECN slope, the ECN field of the packet is unchanged.
 - If the current queue depth is above the configured **max-threshold-percent** of the ECN slope, the ECN field of the packet is (re)marked as Congestion Experienced (CE), ECN=11.
 - If the current queue size is between the **min-threshold-percent** and **max-threshold-percent** of the ECN slope, the ECN field of the packet is (re)marked as CE, ECN=11, based on a probability function that increases linearly from 0% at the minimum threshold to $n\%$ at the maximum threshold, where n is the operational **max-probability** of marking the packet.



Note: The operational values of the **max-probability** may be significantly different from the configured values based on internal hardware calculations. You can check the hardware configured values for any slope calculations.

8. When the packet is non-ECN-capable (the ECN field is zero) and the egress queue has an active WRED slope for the drop probability of the packet, the packet may be dropped by the WRED algorithm, as follows:
- If the current queue depth is below the configured **min-threshold-percent** of the WRED slope, the packet is admitted to the queue.
 - If the current queue depth is above the configured **max-threshold-percent** of the WRED slope, the packet is dropped.
 - If the current queue size is between the minimum threshold and maximum threshold of the WRED slope, the packet is dropped based on a probability function that increases linearly from 0% at the minimum threshold to $n\%$ at the maximum threshold, where n is the operational **max-probability** of dropping the packet.



Note: The operational values of the **max-probability** may be significantly different from the configured values based on internal hardware calculations. You can check the hardware configured values for any WRED slope calculations.

9. Each unicast queue and each multicast queue of an egress port is associated with a scheduler node. The mapping of queues to scheduler nodes is platform-dependent and cannot be configured. See [Output queue scheduling](#).
10. You can configure each egress queue individually with a Peak Information Rate (PIR). The PIR is configured as a percentage of the egress port bandwidth.

By default, the PIR of each queue is 100%. The operational PIR is stored by the **peak-rate-bps** leaf in bits per second. The bits counted in this rate include the Layer 2 framing of the packet (including the 14 byte Ethernet header, the 4-byte VLAN header, and the 4-byte CRC) but exclude the 20-byte Layer 1 overhead (SFD, preamble, IPG).

11. The DSCP field in the IPv4/IPv6 header of the outgoing packet can be rewritten. On the 7250 IXR, the DSCP field must be rewritten when ECN is enabled and the packet ECN field is non-zero. When there is a rewrite policy applied, the DSCP in the outgoing packet is based on the FC (and potentially also the drop probability) of the packet. If the FC (and drop probability) matches an entry in the applied policy, the new DSCP value is based on the policy entry. If there is no matching entry in the applied policy, the new DSCP value is 0.

System default DSCP classifier policy

Table 1: System default DSCP classifier policy

DSCP values	Included DSCP names	Forwarding class	Drop probability
0, 2 to 7	CS0/BE	fc0	Low
1	LE	fc0	High
8 to 11	CS1, AF11	fc1	Low
12 to 13	AF12	fc1	Medium
14 to 15	AF13	fc1	High
16 to 19	CS2, AF21	fc2	Low
20 to 21	AF22	fc2	Medium
22 to 23	AF23	fc2	High
24 to 27	CS3, AF31	fc3	Low
28 to 29	AF32	fc3	Medium
30 to 31	AF33	fc3	High
32 to 35	CS4, AF41	fc4	Low
36 to 37	AF42	fc4	Medium
38 to 39	AF43	fc4	High
40 to 47	CS5, EF	fc5	Low
48 to 55	CS6/NC1	fc6	Low
56 to 63	CS7/NC2	fc7	Low

2.2 How QoS works for VXLAN traffic

When the 7220 IXR-D2 and D3 receives a terminating VXLAN packet on a subinterface, it classifies the packet to one of eight forwarding classes and one of three drop probabilities (low, medium, or high). The classification is based on the following considerations:

- The outer IP header DSCP is ignored.
- If the payload packet is non-IP, the classified FC is fc0 and the classified drop probability is low.

- If the payload packet is IP, and there is a classifier policy referenced by the **qos classifiers vxlan-default** command, that policy determines the FC and drop probability from the header fields of the payload packet.
- If the payload packet is IP, and there is no classifier policy referenced by the **qos classifiers vxlan-default** command, the default DSCP classifier policy determines the FC and drop probability from the header fields of the payload packet.

When the 7220 IXR-D2 and D3 adds VXLAN encapsulation to a packet and forwards it out a subinterface, the inner header IP DSCP value is not modified if the payload packet is IP, even if the egress routed subinterface has a DSCP rewrite rule policy bound to it that matches the packet FC and drop probability. The outer header IP DSCP is modified by the DSCP rewrite rule policy that is bound to the egress routed subinterface, if such a policy exists.

2.3 How QoS works for router-terminated traffic

This section describes how QoS applies to traffic that terminates on the SR Linux.

1. A packet is received on a subinterface and is determined to need extraction toward the CPM. The packet is directed to one of the queues associated with the CPM as a destination “physical port” based on its protocol and type. The following traffic types have their own independent queue:
 - sflow
 - ICMPv4 ping
 - BFD
 - ARP
 - ICMPv6 neighbour solicitation and neighbor advertisement
 - BGP
 - gRPC
 - LLDP
 - IPv4 packets with IP options and IPv6 packets with extension headers
 - DHCPv6
 - IS-IS hello PDUs
 - OSPF/OSPFv3 hello PDUs
2. Some of the queues toward the CPM have a PIR shaping rate designed to prevent an overload of one type of traffic. The PIR shaping rates vary by platform.

2.4 How QoS works for router-originated traffic

This section describes how QoS applies to traffic that originates on the SR Linux.

1. An application on the SR Linux CPM has an IPv4 or IPv6 packet to send to another system.

2. The CPM datapath assigns a DSCP to the self-generated packet based on its protocol and the hard-coded mapping shown in [Table 2: Default forwarding class and DSCP marking for router-originated traffic](#).
 Except for ICMP and ICMPv6 echo-request packets, the DSCP values cannot be overridden. For originated echo-request packets, the DSCP override value can be configured as an optional parameter of the **ping** command.
3. The CPM datapath looks up the DSCP from the previous step (either the fixed value or the override value for echo-request) in the default DSCP classifier policy (see [Table 1: System default DSCP classifier policy](#)) to determine the FC and drop probability level.
4. A forwarding lookup determines the egress port.
5. On the 7250 IXR, the packet is sent to the egress line card and added to a Virtual Output Queue (VOQ) appropriate for its forwarding class and the egress port. The decision to drop or enqueue the packet in the VOQ and the scheduling of the VOQ follows the previous description for transit traffic. There is no scheduling differentiation between router-originated traffic and transit traffic of the same FC on the egress IMM.
6. The packet is directed to the egress queue appropriate for its forwarding class and packet type. On the 7220 IXR-D2 and D3 and the 7220 IXR-H2 and H3, the decision to drop or enqueue the packet in the egress queue and the scheduling of the egress queue follows QoS treatment of transit traffic described in [How QoS works for transit traffic](#).
7. The DSCP field in the IPv4 or IPv6 header is always written based on the hard-coded mapping described in [Table 2: Default forwarding class and DSCP marking for router-originated traffic](#). If the packet also matches a dscp-policy rewrite-rule applied to the output subinterface, the rewrite-rule policy is ignored.

Default forwarding class and DSCP marking for router-originated traffic

Table 2: Default forwarding class and DSCP marking for router-originated traffic

Protocol / message type	Forwarding class	Drop probability	DSCP marking
IPv4 ARP request/reply	6	Low	N/A
ICMPv4 including echo-request ¹ , echo-reply ² , dest-unreachable, redirect, time-exceeded, parameter-problem	0	Medium	0
ICMPv4 echo-request with ToS/DSCP override = X	look up X in system-default DSCP classifier	look up X in system-default DSCP classifier	x
ICMPv4 echo-reply to echo-request with non-zero DSCP x	look up X in system-default DSCP classifier	look up X in system-default DSCP classifier	x
UDP traceroute	0	Low	0
IPv6 neighbor solicitation	6	Low	48 (CS6/NC1)

¹ Echo-request generated by a **ping** command with no DSCP parameter specified.

² Echo-reply to an echo-request packet with DSCP=0.

Protocol / message type	Forwarding class	Drop probability	DSCP marking
IPv6 neighbor advertisement	6	Low	48 (CS6/NC1)
All other ICMPv6 including dest unreachable, packet-too-big, time-exceeded, parameter-problem, echo-request, echo-reply, router-solicitation, redirect	0	Medium	0
ICMPv6 echo-request with DSCP override = x	look up x in system-default DSCP classifier	look up x in system-default DSCP classifier	x
ICMPv6 echo-reply to echo-request with non-zero DSCP x	look up x in system-default DSCP classifier	look up x in system-default DSCP classifier	x
BFD	6	Low	48 (CS6/NC1)
BGP	6	Low	48 (CS6/NC1)
DNS query	4	Low	34 (AF41)
FTP/TFTP	4	Low	34 (AF41)
gNMI	4	Low	34 (AF41)
JSON RPC	4	Low	34 (AF41)
LLDP	N/A	Low	N/A
NTP	4	Low	34 (AF41)
sFlow	0	Low	0
SNMP	4	Low	34 (AF41)
SSH	4	Low	34 (AF41)
Syslog	4	Low	34 (AF41)
TACACS+	4	Low	34 (AF41)

3 QoS configuration

QoS configuration on SR Linux involves the following tasks:

- [DSCP classifier policies configuration for input traffic](#)
- [DSCP rewrite-rule policies configuration for output traffic](#)
- [Queue templates configuration](#)
- [DSCP classifier policy application to subinterfaces](#)
- [Rewrite-rule policy application to subinterfaces](#)
- [Output queue scheduling](#)

3.1 DSCP classifier policies configuration for input traffic

A DSCP classifier policy that is applied to a subinterface attempts to match the 6-bit DSCP value in the IP header of incoming packets to one of its entries. If there is a match, the system assigns the incoming packet to the specified forwarding class and drop probability. If there is no match, the assigned forwarding class is 0 and the assigned drop probability is low.

Packets that require a similar treatment (per-hop behavior) are grouped into an FC, also known as a behavior aggregate. The SR Linux differentiates up to eight forwarding classes.

The drop probability can be high, medium, or low. The default is low. If a queue-template with different WRED slopes is bound to a queue, when the queue experiences congestion, the queue drops packets in the following order:

- high drop probability packets
- medium drop probability packets
- low drop probability packets

3.1.1 Configuring DSCP classifier policies

The following example creates a DSCP classifier policy:

Example:

```
--{ candidate shared default }--[ ]--
# info qos classifiers
qos {
  classifiers {
    dscp-policy new-policy {
      dscp 0 {
        forwarding-class fc0
        drop-probability high
      }
      dscp 8 {
        forwarding-class fc1
        drop-probability high
      }
    }
  }
}
```

```

    }
  }
}

```



Note: To create a new DSCP classification policy based on the default policy, you can copy the default policy from state in candidate mode, as shown in the following example:

```
# copy from state /qos classifiers dscp-policy default to /qos classifiers dscp-policy test
```

Related topics

[DSCP classifier policy application to subinterfaces](#)

3.1.2 Using a DSCP classifier for VXLAN traffic

On a 7720 IXR-D2 and D3, you can use a classifier policy to classify ingress packets received from any remote VXLAN VTEP. The policy applies to payload packets after VXLAN decapsulation is performed.

The following example shows how the DSCP classifier policy created in the previous example (**new-policy**) can be used for VXLAN traffic:

Example:

```
--{ candidate shared default }--[ ]--
# info qos classifiers
qos {
  classifiers {
    vxlan-default new-policy
  }
}

```

3.2 DSCP rewrite-rule policies configuration for output traffic

When you apply a DSCP rewrite-rule policy to a subinterface, the policy attempts to match the forwarding class (and optionally the drop probability) of outbound packets to one of its entries. If there is a match, the DSCP value of the outbound packet changes to the value specified by the policy. If there is no match, the DSCP value changes to 0.

On 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 systems, if no DSCP rewrite-rule policy is applied to a subinterface, the incoming packet's DSCP remains unchanged at egress.

3.2.1 Configuring DSCP rewrite-rule policies

The following example creates a rewrite-rule policy:

Example:

```
--{ candidate shared default }--[ ]--
# info qos rewrite-rules
qos {
  rewrite-rules {
    dscp-policy normalize {
      map fc0 {

```

```

        dscp 1
    map fc0 {
        dscp 7
    map fc1 {
        dscp 10
        drop-probability low {
            dscp 11
        }
        drop-probability high {
            dscp 13
    map fc2 {
        dscp 23
    map fc3 {
        dscp 31
    }
    }
}

```

Related topics

[Rewrite-rule policy application to subinterfaces](#)

3.3 Queue templates configuration

Queue templates are groups of configuration information that apply to a set of queues. On 7250 IXR systems, the controlled set of queues are VOQs; on 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 systems, the controlled set of queues are egress queues.

The maximum number of queue templates per system varies by platform. On 7250 IXR systems, the maximum is eight queue templates; on 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 systems, the maximum is 62 queue templates.

The following parameters are configurable inside a queue template:

- The MBS of each queue; this essentially defines the length of each queue. When the queue builds to the MBS level, further packets are dropped. Be aware that discards may occur before the queue reaches MBS (for example, resulting from shared buffer exhaustion, or from the effects of WRED slopes defined for the queue).
- WRED slopes that define probability curves for discarding packets as a function of (weighted) average queue depth. WRED slopes are not supported for multicast queues.
- ECN slopes that define probability curves for marking ECN-capable packets as having experienced congestion, instead of discarding them. ECN slopes are not supported for multicast queues.

If a queue (VOQ or egress queue) does not have a queue template binding, it inherits the settings of the default queue template. The default queue template has a platform-specific MBS default value, no defined queue utilization thresholds, no WRED slopes, and no ECN slopes. You cannot display the default queue template, but its effect is visible by reading the state of individual queues that lack a queue template binding.

3.3.1 Configuring queue templates

The following example creates a queue template that you could use for any of the following:

- a set of VOQs on a 7250 IXR
- an egress queue on a 7220 IXR-D2 and D3

- an egress queue on a 7220 IXR-H2 and H3

Example:

```
--{ candidate shared default }--[ ]--
# info qos
qos {
    queue-templates {
        queue-template wred-ecn-1 {
        }
    }
}
```



Note: This example is only the starting point of a full configuration. Subsequent sections build on this example to create a full configuration.

3.3.2 Queue depth (maximum burst size)

In a queue-template, the **maximum-burst-size** parameter sets the maximum length of an egress queue or set of VOQs. The queue depth is also known as the Maximum Burst Size (MBS). You must set the **maximum-burst-size** parameter to a non-zero value to configure WRED slope and ECN slope parameters.

On the 7250 IXR, the **maximum-burst-size** parameter applies to a set of VOQs. If the parameter is not configured, or is set to 0, the effective MBS of these VOQs is 256MB.

On the 7220 IXR-D2 and D3 or the 7220 IXR-H2 and H3, the **maximum-burst-size** parameter applies to a set of egress queues. If the parameter is not configured or is set to 0, the effective MBS of these egress queues is calculated based on a fair allocation algorithm. You can assign a non-zero MBS value to multicast queues, but Nokia does not recommend this configuration (especially if multicast traffic is being shaped by configuring **peak-rate-percent**), because it can lead to a shortage of multicast-related buffering resources on 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 systems.

3.3.2.1 Configuring queue depth (maximum burst size)

The following example specifies the queue depth with a set **maximum-burst-size**:

Example:

```
--{ candidate shared default }--[ ]--
# info qos
qos {
    queue-templates {
        queue-template wred-ecn-1 {
            queue-depth {
                maximum-burst-size 20
            }
        }
    }
}
```


3.3.3 WRED slope

In a queue template, you can configure WRED policies to handle congestion when queue space is depleted. Without WRED, when a queue reaches its maximum fill size, the queue discards any packets arriving at the queue (known as tail drop).

WRED policies manage queue depth. They help to prevent congestion by starting random discards when the queue reaches a user-configured threshold value. This avoids the impact of discarding all the new incoming packets. By starting random discards at this threshold, an end-system can adjust its sending rate to the available bandwidth.

The WRED curve algorithm is based on two user-configurable thresholds (**min-threshold-percent** and **max-threshold-percent**) and a discard probability factor (**max-probability**).

On the 7220 IXR-D2 and D3 or the 7220 IXR-H2 and H3, you can configure a WRED slope to apply only to TCP or to non-TCP traffic. This can be useful because TCP has built-in mechanisms to adjust its sending rate in response to packet drops. TCP-based senders lower the packet transmission rate when some of the packets fail to reach the far end.

3.3.3.1 Configuring a WRED slope (7250 IXR)

The following example specifies a WRED slope for low drop probability traffic flowing through a set of VOQs on a 7250 IXR. This WRED slope applies to both TCP and non-TCP traffic.

Example:

```
--{ * candidate shared default }--[ ]--
# info qos
  qos {
    queue-templates {
      queue-template wred-ecn-1 {
        active-queue-management {
          wred-slope all drop-probability low {
            min-threshold-percent 10
            max-threshold-percent 25
            max-probability 50
          }
        }
      }
    }
  }
}
```

3.3.3.2 Configuring a WRED slope (7220 IXR)

The following example specifies a WRED slope for TCP traffic that is classified as low drop probability flowing through an egress queue on the 7220 IXR-D2 and D3 or the 7220 IXR-H2 and H3.

Example:

```
--{ * candidate shared default }--[ ]--
# info qos
  qos {
    queue-templates {
      queue-template wred-ecn-1 {
        active-queue-management {
          wred-slope tcp drop-probability low {
            min-threshold-percent 10
            max-threshold-percent 25
          }
        }
      }
    }
  }
}
```

```
max-probability 50
}
}
}
}
```

3.3.4 ECN slope

Some IP applications support the ECN mechanism. With ECN, IP packets originated by such applications are not discarded when they enter a congested queue; instead, they are marked in a special way. The marking uses the two ECN bits in the traffic class field of the IPv4 or IPv6 packet header. The receiver of IP packets marked as having experienced congestion can signal to the sender (through Layer 4 or higher protocols) that it should reduce its sending rate. The advantage of this feedback mechanism is that the sending rate can drop more gradually than the normal response of a TCP sender to packet discards. A more gradual back-off can result in higher effective throughput in the network.

An ECN slope is similar to a WRED slope. It is based on two user-configurable thresholds (**min-threshold-percent** and **max-threshold-percent**) and a marking probability factor (**max-probability**).

To use an ECN slope, you must configure **explicit-congestion-notification**.

3.3.4.1 Configuring an ECN slope (7250 IXR)

On 7250 IXR systems, the configuration requires you to specify an ECN DSCP policy; this is the DSCP rewrite policy that is used when an ECN field rewrite must be performed. In addition, you can only have one ECN slope per queue and it applies to all drop-probability levels.

Example:

The following example specifies an ECN slope applicable to a 7250 IXR system:

```
--{ candidate shared default }--[ ]--
# info qos
qos {
    explicit-congestion-notification {
        ecn-dscp-policy normalize
    }
    queue-templates{
        queue-template wred-ecn-1 {
            queue-depth{
                maximum-burst-size 20{
            }active-queue-management{
                ecn-slope{
                    ecn-drop-probability all{
                    ecn-min-threshold-percent 50
                    ecn-max-threshold-percent 50
                    max-probability 100{
            }
        }
    }
}
}
```

3.3.4.2 Configuring an ECN slope (7220 IXR)

On the 7220 IXR-D2 and D3 or the 7220 IXR-H2 and H3, you can have one ECN slope per drop-probability level of traffic flowing through an egress queue.

Example:

The following example specifies an ECN slope applicable to a 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 system:

```
--{ candidate shared default }--[ ]--
# info qos
  qos {
    explicit-congestion-notification {
    }
    queue-templates {
      queue-template 2 {
        queue-depth {
          maximum-burst-size 100
        }
        active-queue-management {
          ecn-slope high {
            min-threshold-percent 0
            max-threshold-percent 80
            max-probability 90
          }
        }
      }
    }
  }
}
```

3.3.5 Queue utilization thresholds

When a router receives a burst of traffic, and the incoming rate exceeds the available transmission rate, the router queues the excess traffic. If the burst lasts long enough, or it is followed by additional bursts, the queues may overflow, resulting in traffic loss.

To respond to onsets of congestion, you can subscribe to telemetry information that generates an event when specific queues exceed a specified occupancy level.

To assign a utilization threshold to a queue, you must apply a non-default queue template to the queue, and that queue template must specify a non-zero **high-threshold-bytes** value. When the utilization of the queue crosses the specified **high-threshold-bytes** value, a hardware interrupt is raised. XDP records the current system-time and clears the interrupt. In a scaled setup, XDP may take 10 to 15 ms to process and clear each interrupt, meaning multiple threshold crossings within a very short period of time across one or more queues using the same queue template may appear as only a single event in the telemetry stream. When the **high-threshold-bytes** value is 0, the functionality is disabled and no threshold events are generated for the queues covered by the queue template.

SR Linux supports queue utilization thresholds on 7250 IXR, 7220 IXR-D2 and D3, and 7220 IXR-H2 and H3 systems; however, the behavior varies by system.



Note: You can only configure queue utilization thresholds for unicast queues; multicast queues do not support queue utilization thresholds.

3.3.5.1 Configuring queue utilization thresholds on 7250 IXR systems

On a 7250 IXR system, binding a queue template with a non-zero **high-threshold-bytes** value to an egress queue assigns that threshold value to all the VOQs that logically feed this egress queue.

You can configure each queue template that the system supports with a different **high-threshold-bytes** value as needed.

Example:

The following example configures the **high-threshold-bytes** value to 256255:

```
--{ candidate shared default }--[ ]--
# qos queue-templates queue-template 2 queue-depth high-threshold-bytes 256255
--{ candidate shared default }--[ ]--
# commit stay
All changes have been committed. Starting new transaction.
```

Each configured threshold value is rounded up to the nearest multiple of 256 bytes, up to a maximum capped value of MBS. You can observe the rounding (on a per VOQ-set basis) using the **info from state interface queue-statistics unicast-queue virtual-output-queue queue-depth** output. (A VOQ-set consists of the VOQ for core 0 and the VOQ for core 1.)

Example:

In the following example, the **high-threshold-bytes** value was configured to 256255, but is rounded to the lower 256000 value (that is, a multiple of 256 bytes):

```
--{ candidate shared default }--[ ]--
# info from state interface ethernet-2/1 queue-statistics unicast-queue 0 virtual-output-queue
queue-depth
  interface ethernet-2/1 {
    queue-statistics {
      unicast-queue 0 {
        virtual-output-queue {
          queue-depth {
            maximum-burst-size 1203200768
            high-threshold-bytes 256000
          }
        }
      }
    }
  }
}
```

The state tree maintains the time of the last threshold crossing in the **interface queue-statistics unicast-queue virtual-output-queue queue-depth** leaf. This represents the last time when either VOQ in the VOQ-set (core0/core1) exceeded the operational threshold. The value of this leaf is not cleared when you delete or modify the queue template that is bound to the queue/VOQs or the **high-threshold-bytes** configuration in the applied queue template.

3.3.5.2 Configuring queue utilization thresholds on 7220 IXR-D2 and D3 systems

On 7220 IXR-D2 and D3 systems, binding a queue template with a non-zero **high-threshold-bytes** value to an egress queue causes that threshold value to be used for that specific queue, as long as it is a unicast queue. The configuration of this leaf is ignored when this queue template is attached to a multicast queue.

No more than seven different configured **high-threshold-bytes** values are allowed across all the queue templates used. The management server rejects a commit that would leave more than seven different values after all adds, deletes, and modifies are processed.

Example:

The following example configures the **high-threshold-bytes** value to 2048999:

```
--{ candidate shared default }--[ ]--
A# qos queue-templates queue-template 2 queue-depth maximum-burst-size 2049024
high-threshold-bytes 2048999
```

```
--{ candidate shared default }--[ ]--
# commit stay
All changes have been committed. Starting new transaction.
```

Each configured threshold value (that the management server accepts) is rounded up to the nearest multiple of 2048 bytes, up to a maximum capped value of MBS. For this reason, do not configure values that round to the same multiple of 2048 bytes. This causes duplication among the **high-threshold-bytes** values, of which only seven are allowed. You can display the effect of this rounding using the **info from state interface qos output unicast-queue queue-depth** command.

Example:

In the following example, the **high-threshold-bytes** value was configured to *2048999*, but is rounded to a lower *2048000* value (that is, a multiple of 2048 bytes):

```
--{ candidate shared default }--[ ]--
A:# info from state interface ethernet-1/3 qos output unicast-queue 0 queue-
depth
  interface ethernet-1/3 {
    qos {
      output {
        unicast-queue 0 {
          queue-depth {
            maximum-burst-size 2049024
            high-threshold-bytes 2048000
          }
        }
      }
    }
  }
}
```

The state tree maintains the time of the last threshold crossing in the **interface qos output unicast-queue queue-depth last-high-threshold-time** leaf. This represents the last time the queue exceeded the operational threshold. The value of this leaf is not cleared when you delete or modify the queue template that is bound to the queue or the **high-threshold-bytes** configuration in the applied queue-template.

3.3.5.3 Configuring queue utilization thresholds on 7220 IXR-H2 and H3 systems

On 7220 IXR-H2 and H3 systems, binding a queue template with a non-zero **high-threshold-bytes** value to an egress queue causes that threshold value to be used by each ITM that serves the queue. For a high-threshold event, the queue utilization threshold must be exceeded on either ITM.

No more than seven different configured **high-threshold-bytes** values are allowed across all the queue templates used. The management server rejects a commit that would leave more than seven different values after all adds, deletes, and modifies are processed.

Example:

The following example configures the **high-threshold-bytes** value to 254255:

```
--{ candidate shared default }--[ ]--
A# qos queue-templates queue-template 2 queue-depth maximum-burst-size 2049024
  high-threshold-bytes 254255
--{ candidate shared default }--[ ]--
# commit stay
All changes have been committed. Starting new transaction.
```

Each configured threshold value (that the management server accepts) is rounded up to the nearest multiple of 254 bytes, up to a maximum capped value of MBS. For this reason, do not configure values

that round to the same multiple of 254 bytes. This causes duplication among the **high-threshold-bytes** values, of which only seven are allowed. You can display the effect of this rounding using the **info from state interface qos output unicast-queue queue-depth** command.

Example:

In the following example, the **high-threshold-bytes** value was configured to 254255, but is rounded to a lower 254000 value (that is, a multiple of 254 bytes):

```
--{ candidate shared default }--[ ]--
A:# info from state interface ethernet-1/3 qos output unicast-queue 0 queue-
depth
  interface ethernet-1/3 {
    qos {
      output {
        unicast-queue 0 {
          queue-depth {
            maximum-burst-size 2049024
            high-threshold-bytes 254000
          }
        }
      }
    }
  }
}
```

The state tree maintains the time of the last threshold crossing in the **interface qos output unicast-queue queue-depth last-high-threshold-time** leaf. This represents the last time when either ITM exceeded the operational threshold. The value of this leaf is not cleared when you modify or delete the queue-template that is bound to the queue or the **high-threshold-bytes** configuration in the applied queue-template.

3.4 DSCP classifier policy application to subinterfaces

If you apply a DSCP classifier policy to input traffic on a subinterface, incoming packets are evaluated against the policy, and matching packets are assigned to the forwarding class and drop probability specified by the policy. If no classifier policy is applied to the subinterface, the system default DSCP classifier (with the reserved name *default*) is used.

3.4.1 Applying a DSCP classifier policy to input traffic (7250 IXR)

The following example applies a DSCP classifier policy to inbound IPv6 traffic on a subinterface with a 7250 IXR system:

Example:

```
--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
  interface ethernet-1/1 {
    subinterface 1 {
      qos {
        input {
          classifiers {
            ipv6-dscp new-policy
          }
        }
      }
    }
  }
}
```

```
}

```

3.4.2 Applying a DSCP classifier policy to input traffic (7220 IXR)

The following example applies a DSCP classifier policy to inbound traffic on a subinterface with a 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 system:



Note: The 7220 IXR-D2 and D3 and 7220 IXR-H2 and H3 systems do not support separate classifier policies for IPv4 and IPv6 traffic, but you can apply a common policy that applies to both IPv4 and IPv6 traffic.

Example:

```
--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
  interface ethernet-1/1 {
    subinterface 1 {
      qos {
        input {
          classifiers {
            dscp new-policy
          }
        }
      }
    }
  }
}
```

3.5 Rewrite-rule policy application to subinterfaces

When a rewrite-rule policy is applied to output traffic on a subinterface, outbound packets are evaluated against the policy. The policy subjects all packets to remarking, with some exceptions. If no rewrite-rule policy is applied to the subinterface, the DSCP marking of the traffic leaving the subinterface is unchanged, unless it is ECN-capable traffic forwarded by a 7250 IXR system or VXLAN traffic originated by a 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 system. For these exceptions, DSCP may be remarked even in the absence of a rewrite-rule policy applied to the egress subinterface.

On all platforms, rewrite-rule policies do not affect DSCP marking of self-generated traffic.

3.5.1 Applying a rewrite-rule policy to output traffic (7250 IXR)

The following example applies a rewrite-rule policy to outbound traffic on a subinterface with a 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 system:



Note: Common rewrite policies that apply to both IPv4 and IPv6 traffic are supported on 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 systems.

Example:

```
--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
  interface ethernet-1/1 {
    subinterface 1 {
```

```

    qos {
      output {
        rewrite-rules {
          dscp new-rule
        }
      }
    }
  }
}

```

3.5.2 Applying a rewrite-rule policy to output traffic (7220 IXR)

The following example applies a rewrite-rule policy to outbound IPv4 traffic on a subinterface with a 7250 IXR system:



Note: 7250 IXR systems support separate rewrite policies for IPv4 and IPv6 egress traffic.

Example:

```

--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
  interface ethernet-1/1 {
    subinterface 1 {
      qos {
        output {
          rewrite-rules {
            ipv4-dscp new-rule
          }
        }
      }
    }
  }
}

```

3.6 Output queue scheduling

Each unicast queue and each multicast queue of an egress port is associated with a scheduler node. The mapping of queues to scheduler nodes is platform-dependent and cannot be configured.

On 7250 IXR systems, there are two scheduling nodes per port; one for unicast traffic and one for multicast traffic. The two scheduling nodes have a WRR relationship, but the parameters cannot be adjusted. There is one PIR scheduling loop per scheduling node. The scheduling loop serves the strict priority classes first (in descending order of FC), followed by the WRR classes (by weight), limiting each forwarding class to its PIR (expressed as a percentage of the egress port bandwidth). By default, the PIR of each forwarding class is 100%. Note that multicast traffic handled by the multicast scheduler node is unscheduled and is not subject to the ingress VOQ buffering that applies to unicast traffic.

On 7220 IXR-D2 and D3 systems, the unicast queue and multicast queue for a particular forwarding class make up a queue pair. Each of the eight possible queue pairs of an egress port are associated with a scheduler node. Each scheduler node is served as strict priority (SP) or weighted round robin (WRR). If it is served as WRR, the scheduler node also has an associated weight. The scheduling loop serves the SP nodes first, followed by the WRR nodes by weight. The serving order of SP queues is in descending order of FC: fc7 first, then fc6, then fc5, and so on.

On 7220 IXR-H2 and H3 systems, there is a one-to-one mapping of queues to scheduler nodes. Each scheduler node can be served as SP or WRR. A WRR node has a configurable weight. The scheduling loop serves the SP nodes first, followed by the WRR nodes by weight. The serving order of SP queues is as follows:

- unicast queue 7 serving fc7
- unicast queue 6 serving fc6
- multicast queue 3 serving fc6 and fc7
- unicast queue 5 serving fc5
- unicast queue 4 serving fc4
- multicast queue 2 serving fc4 and fc5
- unicast queue 3 serving fc3
- unicast queue 2 serving fc2
- multicast queue 1 serving fc2 and fc3
- unicast queue 1 serving fc1
- unicast queue 0 serving fc0
- multicast queue 0 serving fc0 and fc1

3.6.1 Configuring strict priority (7250 IXR)

The following example configures a queue or scheduler node for strict priority. When strict priority is set to false, the associated queue or scheduler node is configured as WRR. When strict priority is set to true, any configured weight is ignored.

Example:

```
--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
  interface ethernet-1/1 {
    qos {
      output {
        unicast-queue 0 {
          scheduling {
            strict-priority true
          }
        }
      }
    }
  }
}
```

3.6.2 Configuring strict priority (7220 IXR-D2 and D3 or 7220 IXR-H2 and H3)

The following example configures a queue or scheduler node for strict priority. Note that when strict priority is set to false, the associated queue or scheduler node is configured as WRR. When strict priority is set to true, any configured weight is ignored.

Example:

```
--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
```

```

interface ethernet-1/1 {
  qos {
    output {
      scheduler {
        tier 1 {
          node 0 {
            strict-priority true
          }
        }
      }
    }
  }
}

```

3.6.3 Configuring WRR (7250 IXR)

The following example configures a queue or scheduler-node for WRR. Queues or scheduler nodes that you do not configure with a specific weight have a weight of 1.

Example:

```

--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
interface ethernet-1/1 {
  qos {
    output {
      unicast-queue 0 {
        scheduling {
          strict-priority false {
            weight 20
          }
        }
      }
    }
  }
}

```

3.6.4 Configuring WRR (7220 IXR-D2 and D3 or 7220 IXR-H2 and H3)

The following example configures a queue or scheduler node for WRR. Queues or scheduler nodes that you do not configure with a specific weight have a weight of 1.

Example:

```

--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
interface ethernet-1/1 {
  qos {
    output {
      scheduler {
        tier 1 {
          node 0 {
            strict-priority false
            weight 20
          }
        }
      }
    }
  }
}

```

```
}
```

3.6.5 Configuring forwarding class peak rate

The following example sets the maximum percentage of port bandwidth that is available to traffic of a particular FC. By default, traffic belonging to any FC can use up to 100% of the port bandwidth. The example is applicable to 7250 IXR, 7220 IXR-D2 and D3, and 7220 IXR-H2 and H3 system.

Example:

```
--{ candidate shared default }--[ ]--  
# info interface ethernet-1/1  
  interface ethernet-1/1 {  
    qos {  
      output {  
        unicast-queue 0 {  
          scheduling {  
            peak-rate-percent 75  
          }  
        }  
      }  
    }  
  }  
}
```

4 MPLS QoS overview

SR Linux supports QoS capabilities in MPLS networks using traffic classification and marking.

MPLS traffic classification and marking

SR Linux supports EXP-inferred LSPs as described in RFC 3270. This allows multiple classes of service to be transported by a single LSP, with the EXP marking of each packet determining the correct per-hop behavior (PHB) to apply to each router.

On SR Linux, the mapping between an EXP value and a PHB is provided by an MPLS traffic-class classifier policy. A single router can have one or more of these policies so that some subinterfaces can have one policy applied and other subinterfaces can have another policy applied. Each traffic-class classifier policy consists of multiple mapping entries, each of which maps one unique EXP value to a (forwarding class, drop probability) tuple.

SR Linux also supports MPLS traffic-class rewrite policies. If MPLS-encapsulated packets are transmitted out an egress subinterface with such a policy bound to it, the EXP field in all the pushed labels of these packets is based on the mapping rules of the policy. MPLS traffic-class rewrite rules associate a forwarding class or a (forwarding class, drop probability) tuple with an EXP rewrite value.

SR Linux does not support the short-pipe model of RFC 3270.

4.1 Ingress LER

When an SR Linux router that is acting as an ingress LER matches an IP packet to an LDP tunnel or a static MPLS forwarding entry, the following apply.

- The ingress LER determines the forwarding class and drop probability of the packet from the IP DSCP of the received unlabeled packet, based on the DSCP classifier policy applied to the ingress subinterface (or the default DSCP classifier policy if there is no explicit association). If an MPLS TC classifier policy is applied to the ingress subinterface, it has no effect.
- If a DSCP rewrite policy is applied to the egress subinterface, the IP header DSCP value is rewritten before the egress MPLS encapsulation is applied.
- If no MPLS TC rewrite policy is associated with the egress subinterface, EXP=0 is written into all pushed labels.
- If an MPLS TC rewrite policy is associated with the egress subinterface, and it matches the forwarding class (and possibly also the drop probability) of the packet, the EXP provided by the mapping rule is written into the EXP field of all pushed labels.
- If ECN is enabled globally, and the packet hits an ECN slope in a congested queue such that the ECN marking should be '11', the ECN field of the packet is modified accordingly, and the DSCP field is also remarked according to the ECN DSCP policy.

4.2 Transit LSR

When an SR Linux router that is acting as a transit LSR matches an MPLS packet to a swap ILM entry, the following apply.

- The transit LSR determines the forwarding class and drop probability of the packet from the EXP in the topmost label stack entry of the received labeled packet (before popping), based on the MPLS TC classifier policy applied to the ingress subinterface (or the default MPLS TC classifier policy, if there is no explicit association).
- If a DSCP classifier policy is applied to the ingress subinterface, it has no effect on the packet classification.
- If a DSCP rewrite policy is applied to the egress subinterface, it has no effect on the transmitted MPLS packet.
- If no MPLS TC rewrite policy is associated with the egress subinterface, the classified FC of the packet is written as a value 0 to 7 into the EXP field of all pushed labels. This does not guarantee that the EXP of the popped labels matches the EXP of the pushed labels (that is, if a non-default MPLS TC classifier policy is applied to the ingress subinterface).
- If an MPLS TC rewrite policy is associated with the egress subinterface, and it matches the forwarding class (and possibly also the drop probability) of the packet, the EXP provided by the mapping rule is written into the EXP field of all pushed labels.
- If ECN is enabled globally, it has no effect on the MPLS packet. The MPLS packet is considered non-ECT capable, even if the buried IP ECN bits indicate otherwise. The IP ECN field is not modified.

4.3 PHP LSR

When an SR Linux router that is acting as a PHP LSR matches an MPLS packet to a pop and swap-to-implicit-null ILM entry, the following apply.

- The PHP LSR determines the forwarding class and drop probability of the packet from the EXP in the topmost label stack entry of the received labeled packet (before popping), based on the MPLS TC classifier policy applied to the ingress subinterface (or the default MPLS TC classifier policy, if there is no explicit association). If a DSCP classifier policy is applied to the ingress subinterface, it has no effect on the classification of the packet.
- If an MPLS TC rewrite policy is applied to the egress subinterface, it has no effect on the transmitted IP packet.
- If no DSCP rewrite policy is associated with the egress subinterface, the DSCP field of the IP payload packet is transmitted unchanged. There is no attempt to copy the EXP field into the IP DSCP of the IP payload packet.
- If a DSCP rewrite policy is associated with the egress subinterface, and it matches the forwarding class (and possibly also the drop probability) of the packet, the DSCP provided by the mapping rule is written (as an override) into the DSCP field in the transmitted IP packet. This is consistent with the uniform model of RFC 3270.
- If ECN is enabled globally, it has no effect on the PHP packet. The PHP packet is considered non-ECT capable even if the IP ECN bits indicate otherwise. The IP ECN field is not modified.

4.4 Egress LER

When an SR Linux router that is acting as an egress LER matches an MPLS packet to a pop ILM entry that leads to all labels being popped, the following apply.

- The egress LER determines the forwarding class and drop probability of the packet from the EXP in the topmost label stack entry of the received labeled packet (before popping), based on the mpls-tc classifier policy applied to the ingress subinterface (or the default mpls-tc classifier policy, if there is no explicit association).
If a DSCP classifier policy is applied to the ingress subinterface, it has no effect on the classification of the packet.
- If an mpls-tc rewrite policy is applied to the egress subinterface, it has no effect on the transmitted IP packet.
- If no DSCP rewrite policy is associated with the egress subinterface, the DSCP field of the IP payload packet is transmitted unchanged. There is no attempt to copy the EXP field into the IP DSCP of the IP payload packet. This is consistent with the pipe model of RFC 3270.
- If a DSCP rewrite policy is associated with the egress subinterface, and it matches the forwarding class (and possibly also the drop probability) of the packet, the DSCP provided by the mapping rule is copied into the IP DSCP of the transmitted IP packet, overwriting the previous value. This is consistent with the uniform model of RFC 3270.
- If ECN is enabled globally, it has no effect on the terminating MPLS packet. The terminating packet is considered non-ECT capable even if the IP ECN bits indicate otherwise. The IP ECN field is not modified.



Note: Note that the DSCP marking of terminating MPLS traffic cannot be decoupled from the DSCP marking of transit IP traffic through the same egress subinterface.

4.5 Default MPLS traffic-class classifier policy

The following table shows the default MPLS TC classifier policy.

Table 3: Default MPLS TC classifier policy

Traffic class (EXP)	Forwarding class	Drop probability
0	0	low
1	1	low
2	2	low
3	3	low
4	4	low
5	5	low
6	6	low
7	7	low

5 MPLS QoS configuration

MPLS QoS configuration on SR Linux involves the following tasks:

- [Configuring MPLS traffic-class policy](#)
- [Applying MPLS traffic-class policy to input traffic](#)
- [Configuring MPLS rewrite rules](#)
- [Applying MPLS rewrite rules to output traffic](#)

5.1 Configuring MPLS traffic-class policy

The following example creates an MPLS traffic-class policy:

Example:

```
--{ candidate shared default }--[ ]--
# info qos classifiers
  qos {
    classifiers {
      mpls-traffic-class-policy mpls-policy-1 {
        traffic-class 7 {
          forwarding-class fc7
          drop-probability medium
        }
      }
    }
  }
}
```

5.2 Applying MPLS traffic-class policy to input traffic

The following example applies an MPLS traffic-class policy to inbound traffic on a subinterface.

Example:

```
--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
  interface ethernet-1/1 {
    subinterface 1 {
      qos {
        input {
          classifiers {
            mpls-traffic-class mpls-policy-1
          }
        }
      }
    }
  }
}
```

5.3 Configuring MPLS rewrite rules

The following example creates an MPLS rewrite-rule policy:

Example:

```
--{ candidate shared default }--[ ]--
# info qos rewrite-rules
  qos {
    rewrite-rules {
      mpls-traffic-class-policy mpls-rewrite-2 {
        map fc7 {
          traffic-class 7
        }
      }
    }
  }
}
```

5.4 Applying MPLS rewrite rules to output traffic

The following example applies a rewrite-rule policy to outbound traffic on a subinterface.

Example:

```
--{ candidate shared default }--[ ]--
# info interface ethernet-1/1
  interface ethernet-1/1 {
    subinterface 1 {
      qos {
        output {
          rewrite-rules {
            mpls-traffic-class mpls-rewrite-2
          }
        }
      }
    }
  }
}
```


6 Buffer utilization display

The following table describes the buffer utilization differences between the 7250 IXR, 7220 IXR-D2 and D3, or 7220 IXR-H2 and H3.

Table 4: Buffer utilization

Hardware	Buffer memory
7250 IXR	<ul style="list-style-type: none"> • SRAM size = 32MB • DRAM (HBM) size = 8 GB
7220 IXR-D2 and D3	<ul style="list-style-type: none"> • Total Buffer size = 32MB • Reserved Buffer size = 4.65MB
7220 IXR-H2 and H3	<ul style="list-style-type: none"> • Total Buffer size = 64MB • Reserved Buffer size = 6.7MB

6.1 Displaying buffer utilization

The following examples show overall buffer usage. The output varies depending on the hardware deployed.

Example: Displaying buffer utilization (7250 IXR)

```
# info from state platform linecard 1 forwarding-complex 0 buffer-memory
platform {
  linecard 1 {
    forwarding-complex 0 {
      buffer-memory {
        sram {
          used 15808512 >> in bytes
          free 17745920 >> in bytes
        }
        dram {
          used 48 >>> it is in % of DRAM
        }
      }
    }
  }
}
```

Example: Displaying buffer utilization (7220 IXR-D2 and D3 or 7220 IXR-H2 and H3)

```
# info from state platform linecard 1 forwarding-complex 0 buffer-memory
platform {
  linecard 1 {
    forwarding-complex 0 {
      buffer-memory {
        used 2097152
      }
    }
  }
}
```

```
    free 27263246
    reserved 4194034
  }
}
}
```

7 Displaying QoS statistics

To display traffic statistics for each output queue on an interface, use the **show interface <id> queue-detail** command in running or candidate mode.

The following example displays output queue statistics for an interface on a 7250 IXR system. The output on a 7220 IXR-D2 and D3 or 7220 IXR-H2 and H3 system is similar, but shows slightly different information.

Example:

```
# show interface ethernet-1/1 queue-detail
=====
Interface: ethernet-1/1
-----
Description      : <None>
Oper state       : up
Last change      : 17h3m43s ago, 1 flaps since last clear
Auto-negotiation: false
Duplex           : N/A
Speed            : 100G
Loopback mode    : false
MTU              : 9232
VLAN tagging     : false
MAC address      : 12:12:02:FF:00:00
Last stats clear: never
=====
Scheduler details for for ethernet-1/1
-----
Tier 1
Node  Scheduling  Weight  Serving
0     SP          -       unicast queue 0, 1, 2, 3, 4, 5, 6, 7 multicast
      SP          -       queue 0, 1, 2, 3, 4, 5, 6, 7
1     SP          -       -
=====
Unicast Queue   : 0
Forwarding class: fc0
Queue template  : default
-----
Scheduling
-----
PIR (%)         : 100
PIR (bps)       : -
Strict Priority: true
Weight         : 1
Scheduler node  : Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Active WRED Slopes
-----
Slope  Traffic type  Drop probability  Min-threshol  Max-threshol  Max
1      all         low               d(%MBS)       d(%MBS)       probability
2      all         medium
3      all         high
```

```

Active ECN Slopes
-----
Slope  Drop probability      Min-      Max-      Max
      probability      threshold(%MBS)  threshold(%MBS)  probability
-----
1      all
-----

Queue Statistics
-----
Tx Packets      : 6
Tx Bytes        : 580
Dropped Packets: 0
Dropped Bytes   : 0
-----

V0Q Id          Fwd-      Fwd-Pkts(L/M/H)      Drop-      Drop-
                Octets(L/M/H)          Pkts(L/M/H)          Pkts(L/M/H)
-----
1              0/0/0              0/0/0              0/0/0              0/0/0
2              0/0/0              0/0/0              0/0/0              0/0/0
3              0/0/0              0/0/0              0/0/0              0/0/0
4              0/0/0              0/0/0              0/0/0              0/0/0
5              0/0/0              0/0/0              0/0/0              0/0/0
6              0/0/0              0/0/0              0/0/0              0/0/0
7              0/0/0              0/0/0              0/0/0              0/0/0
8              0/0/0              0/0/0              0/0/0              0/0/0
=====

Unicast Queue   : 1
Forwarding class: fc1
Queue template  : default
-----

Scheduling
-----
PIR (%)         : 100
PIR (bps)       : -
Strict Priority: true
Weight          : 1
Scheduler node  : Tier 1, Node 0
-----

Queue Depth
-----
Maximum burst (bytes): 0
-----

Active WRED Slopes
-----
Slope  Traffic type      Drop      Min-threshol  Max-threshol  Max
      Traffic type      probability  d(%MBS)      d(%MBS)      probability
-----
1      all              low
2      all              medium
3      all              high
-----

Active ECN Slopes
-----
Slope  Drop probability      Min-      Max-      Max
      probability      threshold(%MBS)  threshold(%MBS)  probability
-----
1      all
-----

Queue Statistics
-----
Tx Packets      : 0
Tx Bytes        : 0
Dropped Packets: 0
Dropped Bytes   : 0
-----

V0Q Id          Fwd-      Fwd-Pkts(L/M/H)      Drop-      Drop-
                Octets(L/M/H)          Pkts(L/M/H)          Pkts(L/M/H)
-----
1              0/0/0              0/0/0              0/0/0              0/0/0
2              0/0/0              0/0/0              0/0/0              0/0/0
3              0/0/0              0/0/0              0/0/0              0/0/0
    
```

```

4      0/0/0      0/0/0      0/0/0      0/0/0
5      0/0/0      0/0/0      0/0/0      0/0/0
6      0/0/0      0/0/0      0/0/0      0/0/0
7      0/0/0      0/0/0      0/0/0      0/0/0
8      0/0/0      0/0/0      0/0/0      0/0/0
=====
Unicast Queue   : 2
Forwarding class: fc2
Queue template  : default
-----
Scheduling
-----
PIR (%)        : 100
PIR (bps)      : -
Strict Priority: true
Weight         : 1
Scheduler node : Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Active WRED Slopes
-----
Slope  Traffic type      Drop      Min-threshol  Max-threshol  Max
      probability      d(%MBS)   d(%MBS)      probability
1      all              low
2      all              medium
3      all              high
-----
Active ECN Slopes
-----
Slope  Drop probability  Min-      Max-      Max
      probability      threshold(%MBS)  threshold(%MBS)  probability
1      all
-----
Queue Statistics
-----
Tx Packets      : 0
Tx Bytes        : 0
Dropped Packets: 0
Dropped Bytes   : 0
-----
V0Q Id      Fwd-      Fwd-Pkts(L/M/H)      Drop-      Drop-
      Octets(L/M/H)      Octets(L/M/H)      Pkts(L/M/H)
1      0/0/0      0/0/0      0/0/0      0/0/0
2      0/0/0      0/0/0      0/0/0      0/0/0
3      0/0/0      0/0/0      0/0/0      0/0/0
4      0/0/0      0/0/0      0/0/0      0/0/0
5      0/0/0      0/0/0      0/0/0      0/0/0
6      0/0/0      0/0/0      0/0/0      0/0/0
7      0/0/0      0/0/0      0/0/0      0/0/0
8      0/0/0      0/0/0      0/0/0      0/0/0
=====
Unicast Queue   : 3
Forwarding class: fc3
Queue template  : default
-----
Scheduling
-----
PIR (%)        : 100
PIR (bps)      : -
Strict Priority: true
Weight         : 1
Scheduler node : Tier 1, Node 0
-----

```

```

Queue Depth
-----
Maximum burst (bytes): 0
-----
Active WRED Slopes
-----
Slope  Traffic type      Drop
      probability    Min-threshol    Max-threshol    Max
      probability    d(%MBS)         d(%MBS)         probability
1      all            low
2      all            medium
3      all            high
-----
Active ECN Slopes
-----
Slope  Drop probability    Min-
      threshold(%MBS)  Max-
      threshold(%MBS)  Max probability
1      all
-----
Queue Statistics
-----
Tx Packets      : 0
Tx Bytes        : 0
Dropped Packets: 0
Dropped Bytes   : 0
-----
VQ Id   Fwd-      Fwd-Pkts(L/M/H)    Drop-      Drop-
      Octets(L/M/H)    Octets(L/M/H)      Pkts(L/M/H)
1      0/0/0      0/0/0              0/0/0      0/0/0
2      0/0/0      0/0/0              0/0/0      0/0/0
3      0/0/0      0/0/0              0/0/0      0/0/0
4      0/0/0      0/0/0              0/0/0      0/0/0
5      0/0/0      0/0/0              0/0/0      0/0/0
6      0/0/0      0/0/0              0/0/0      0/0/0
7      0/0/0      0/0/0              0/0/0      0/0/0
8      0/0/0      0/0/0              0/0/0      0/0/0
=====
Unicast Queue   : 4
Forwarding class: fc4
Queue template  : default
-----
Scheduling
-----
PIR (%)         : 100
PIR (bps)       : -
Strict Priority: true
Weight          : 1
Scheduler node  : Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Active WRED Slopes
-----
Slope  Traffic type      Drop
      probability    Min-threshol    Max-threshol    Max
      probability    d(%MBS)         d(%MBS)         probability
1      all            low
2      all            medium
3      all            high
-----
Active ECN Slopes
-----
Slope  Drop probability    Min-
      threshold(%MBS)  Max-
      threshold(%MBS)  Max probability
1      all
-----

```

Queue Statistics

Tx Packets : 104
 Tx Bytes : 20266
 Dropped Packets: 0
 Dropped Bytes : 0

VOQ Id	Fwd-Octets(L/M/H)	Fwd-Pkts(L/M/H)	Drop-Octets(L/M/H)	Drop-Pkts(L/M/H)
1	0/0/0	0/0/0	0/0/0	0/0/0
2	0/0/0	0/0/0	0/0/0	0/0/0
3	0/0/0	0/0/0	0/0/0	0/0/0
4	0/0/0	0/0/0	0/0/0	0/0/0
5	0/0/0	0/0/0	0/0/0	0/0/0
6	0/0/0	0/0/0	0/0/0	0/0/0
7	0/0/0	0/0/0	0/0/0	0/0/0
8	0/0/0	0/0/0	0/0/0	0/0/0

Unicast Queue : 5
 Forwarding class: fc5
 Queue template : default

Scheduling

PIR (%) : 100
 PIR (bps) : -
 Strict Priority: true
 Weight : 1
 Scheduler node : Tier 1, Node 0

Queue Depth

Maximum burst (bytes): 0

Active WRED Slopes

Slope	Traffic type	Drop probability	Min-threshold d(%MBS)	Max-threshold d(%MBS)	Max probability
1	all	low			
2	all	medium			
3	all	high			

Active ECN Slopes

Slope	Drop probability	Min-threshold(%MBS)	Max-threshold(%MBS)	Max probability
1	all			

Queue Statistics

Tx Packets : 0
 Tx Bytes : 0
 Dropped Packets: 0
 Dropped Bytes : 0

VOQ Id	Fwd-Octets(L/M/H)	Fwd-Pkts(L/M/H)	Drop-Octets(L/M/H)	Drop-Pkts(L/M/H)
1	0/0/0	0/0/0	0/0/0	0/0/0
2	0/0/0	0/0/0	0/0/0	0/0/0
3	0/0/0	0/0/0	0/0/0	0/0/0
4	0/0/0	0/0/0	0/0/0	0/0/0
5	0/0/0	0/0/0	0/0/0	0/0/0
6	0/0/0	0/0/0	0/0/0	0/0/0
7	0/0/0	0/0/0	0/0/0	0/0/0
8	0/0/0	0/0/0	0/0/0	0/0/0

```

Unicast Queue : 6
Forwarding class: fc6
Queue template : default
-----
Scheduling
-----
PIR (%) : 100
PIR (bps) : -
Strict Priority: true
Weight : 1
Scheduler node : Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Active WRED Slopes
-----
Slope Traffic type Drop probability Min-threshol Max-threshol Max
1 all low d(%MBS) d(%MBS) probability
2 all medium
3 all high
-----
Active ECN Slopes
-----
Slope Drop probability Min- Max- Max probability
1 all threshold(%MBS) threshold(%MBS)
-----
Queue Statistics
-----
Tx Packets : 8249
Tx Bytes : 705888
Dropped Packets: 0
Dropped Bytes : 0
-----
VOQ Id Fwd- Fwd-Pkts(L/M/H) Drop- Drop-
Octets(L/M/H) Pkts(L/M/H)
1 0/0/0 0/0/0 0/0/0 0/0/0
2 0/0/0 0/0/0 0/0/0 0/0/0
3 0/0/0 0/0/0 0/0/0 0/0/0
4 0/0/0 0/0/0 0/0/0 0/0/0
5 0/0/0 0/0/0 0/0/0 0/0/0
6 0/0/0 0/0/0 0/0/0 0/0/0
7 0/0/0 0/0/0 0/0/0 0/0/0
8 0/0/0 0/0/0 0/0/0 0/0/0
=====
Unicast Queue : 7
Forwarding class: fc7
Queue template : default
-----
Scheduling
-----
PIR (%) : 100
PIR (bps) : -
Strict Priority: true
Weight : 1
Scheduler node : Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Active WRED Slopes
-----

```



```

Slope  Traffic type  Drop probability  Min-threshol  Max-threshol  Max
1      all          low              d(%MBS)      d(%MBS)      probability
2      all          medium
3      all          high
-----
Active ECN Slopes
-----
Slope  Drop probability  Min-          Max-          Max
1      all              threshold(%MBS)  threshold(%MBS)  probability
-----
Queue Statistics
-----
Tx Packets      : 0
Tx Bytes        : 0
Dropped Packets: 0
Dropped Bytes   : 0
-----
V0Q Id          Fwd-          Fwd-Pkts(L/M/H)  Drop-          Drop-
                Octets(L/M/H)  Octets(L/M/H)    Octets(L/M/H)  Pkts(L/M/H)
1              0/0/0         0/0/0             0/0/0          0/0/0
2              0/0/0         0/0/0             0/0/0          0/0/0
3              0/0/0         0/0/0             0/0/0          0/0/0
4              0/0/0         0/0/0             0/0/0          0/0/0
5              0/0/0         0/0/0             0/0/0          0/0/0
6              0/0/0         0/0/0             0/0/0          0/0/0
7              0/0/0         0/0/0             0/0/0          0/0/0
8              0/0/0         0/0/0             0/0/0          0/0/0
=====
Multicast Queue : 0
Forwarding class: fc0
-----
Scheduling
-----
PIR (bps)      : -
Scheduler node: Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Queue Statistics
-----
Tx Packets      : 0
Tx Bytes        : 0
Dropped Packets: 0
Dropped Bytes   : 0
-----
=====
Multicast Queue : 1
Forwarding class: fc1
-----
Scheduling
-----
PIR (bps)      : -
Scheduler node: Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Queue Statistics
-----
Tx Packets      : 0
Tx Bytes        : 0

```

```

Dropped Packets: 0
Dropped Bytes : 0
-----
=====
Multicast Queue : 2
Forwarding class: fc2
-----
Scheduling
-----
PIR (bps)      : -
Scheduler node: Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Queue Statistics
-----
Tx Packets     : 0
Tx Bytes       : 0
Dropped Packets: 0
Dropped Bytes  : 0
-----
=====
Multicast Queue : 3
Forwarding class: fc3
-----
Scheduling
-----
PIR (bps)      : -
Scheduler node: Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Queue Statistics
-----
Tx Packets     : 0
Tx Bytes       : 0
Dropped Packets: 0
Dropped Bytes  : 0
-----
=====
Multicast Queue : 4
Forwarding class: fc4
-----
Scheduling
-----
PIR (bps)      : -
Scheduler node: Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Queue Statistics
-----
Tx Packets     : 0
Tx Bytes       : 0
Dropped Packets: 0
Dropped Bytes  : 0
-----
=====
Multicast Queue : 5
Forwarding class: fc5

```

```

-----
Scheduling
-----
PIR (bps)      : -
Scheduler node: Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Queue Statistics
-----
Tx Packets     : 0
Tx Bytes       : 0
Dropped Packets: 0
Dropped Bytes  : 0
-----
=====
Multicast Queue : 6
Forwarding class: fc6
-----
Scheduling
-----
PIR (bps)      : -
Scheduler node: Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Queue Statistics
-----
Tx Packets     : 0
Tx Bytes       : 0
Dropped Packets: 0
Dropped Bytes  : 0
-----
=====
Multicast Queue : 7
Forwarding class: fc7
-----
Scheduling
-----
PIR (bps)      : -
Scheduler node: Tier 1, Node 0
-----
Queue Depth
-----
Maximum burst (bytes): 0
-----
Queue Statistics
-----
Tx Packets     : 0
Tx Bytes       : 0
Dropped Packets: 0
Dropped Bytes  : 0
-----
=====

```

7.1 Clearing QoS statistics

You can reset the queue statistics counters for an interface.

Example: Reset all statistics counters on an interface

The following example resets all statistics counters on an interface:

```
--{ running }--[ ]--
# tools interface ethernet-1/1 statistics queue-statistics clear
```

Example: Reset statistics counters for multicast egress queue

The following example resets statistics counters for a specified egress queue (multicast) on an interface:

```
--{ running }--[ ]--
# tools interface ethernet-1/1 statistics queue-statistics multicast-queue 1 clear
```

7.2 QoS profile resource usage

A QoS profile resource refers to the number of classifier and rewrite policies that are applied to interfaces on a line card. Each classifier or rewrite policy that is applied to an interface on a line card counts as one profile resource used.

For example, if you create classifier policy `ds cp1` and apply it to input IPv4 traffic on an interface, and apply the same `ds cp1` policy to input IPv6 traffic on a different interface on the same line card, it counts as two classifier profile resources used.

The SR Linux supports up to 15 classifier profile resources and up to 32 rewrite profile resources per line card. You can display the number of QoS profile resources in use for each line card.

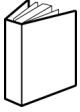
7.2.1 Displaying QoS profile resource usage on a 7250 IXR system

The following example displays the number of used and free classifier and rewrite profile resources for a line card:

Example:

```
# info from state platform linecard 1 forwarding-complex 0 qos
platform {
  linecard 1 {
    forwarding-complex 0 {
      qos {
        resource classifier-profiles {
          used 1
          free 15
        }
        resource rewrite-profiles {
          used 1
          free 31
        }
      }
    }
  }
}
```


Customer document and product support



Customer documentation

[Customer documentation welcome page](#)



Technical support

[Product support portal](#)



Documentation feedback

[Customer documentation feedback](#)