

Pianificazione dell'output QoS sugli switch Catalyst serie 6500/6000 con software CatOS

Sommario

[Introduzione](#)

[Prerequisiti](#)

[Requisiti](#)

[Componenti usati](#)

[Convenzioni](#)

[Premesse](#)

[Rilasci coda di output](#)

[Tipi di accodamento coinvolti nella pianificazione dell'output sugli switch Catalyst 6500/6000](#)

[Caduta coda](#)

[Random Early Detection e Weighted Random Early Detection](#)

[Round Robin ponderato](#)

[Coda con priorità rigorosa](#)

[Capacità di accodamento dell'output di diverse schede di linea su Catalyst 6000](#)

[Funzionalità comando show port](#)

[Comprendere la funzionalità di accodamento di una porta](#)

[Creazione di QoS sugli switch Catalyst 6500/6000](#)

[Meccanismo di programmazione dell'output su Catalyst 6500/6000](#)

[Pianificazione della configurazione, del monitoraggio e dell'output sugli switch Catalyst 6500/6000](#)

[Configurazione predefinita per QoS sugli switch Catalyst 6500/6000](#)

[Configurazione](#)

[Monitorare la pianificazione dell'output e verificare la configurazione](#)

[Utilizzare la programmazione dell'output per ridurre il ritardo e l'instabilità](#)

[Riduci ritardo](#)

[Riduci variazione](#)

[Informazioni correlate](#)

[Introduzione](#)

La pianificazione dell'output garantisce che il traffico importante non venga interrotto in caso di sovrassegnazione. In questo documento vengono descritte tutte le tecniche e gli algoritmi utilizzati per la pianificazione dell'output sugli switch Cisco Catalyst serie 6500/6000 con software di sistema Catalyst OS (CatOS). Questo documento offre anche una breve panoramica sulla funzionalità di accodamento degli switch Catalyst 6500/6000 e su come configurare i diversi parametri della pianificazione dell'output.

Nota: se si esegue il software Cisco IOS® sugli switch Catalyst 6500/6000, per ulteriori informazioni fare riferimento a [QoS Output Scheduling sugli switch Catalyst serie 6500/6000 con](#)

[software Cisco IOS.](#)

Prerequisiti

Requisiti

Nessun requisito specifico previsto per questo documento.

Componenti usati

Gli esempi riportati in questo documento sono stati creati da un Catalyst 6000 con Supervisor Engine 1A e Policy Feature Card (PFC). Gli esempi sono validi anche per un Supervisor Engine 2 con PFC2 o per un Supervisor Engine 720 con PFC3.

Le informazioni discusse in questo documento fanno riferimento a dispositivi usati in uno specifico ambiente di emulazione. Su tutti i dispositivi menzionati nel documento la configurazione è stata ripristinata ai valori predefiniti. Se la rete è operativa, valutare attentamente eventuali conseguenze derivanti dall'uso dei comandi.

Convenzioni

Per ulteriori informazioni sulle convenzioni usate, consultare il documento [Cisco sulle convenzioni nei suggerimenti tecnici.](#)

Premesse

Rilasci coda di output

Le perdite di output sono causate da un'interfaccia congestionata. Una causa comune di questo problema potrebbe essere il traffico proveniente da un collegamento con larghezza di banda elevata che viene convertito in un collegamento con larghezza di banda inferiore o il traffico proveniente da più collegamenti in entrata che vengono convertiti in un unico collegamento in uscita.

Ad esempio, se si invia una grande quantità di traffico bursty su un'interfaccia Gigabit e si passa a un'interfaccia a 100 Mbps, è possibile che si verifichino cali di output incrementali sull'interfaccia a 100 Mbps. Infatti, la coda di output su quell'interfaccia è sovraccarica dal traffico in eccesso a causa della mancata corrispondenza tra la larghezza di banda in entrata e in uscita. La velocità del traffico sull'interfaccia in uscita non può accettare tutti i pacchetti che devono essere inviati.

La soluzione definitiva per risolvere il problema è aumentare la velocità della linea. Esistono tuttavia modi per impedire, ridurre o controllare le riduzioni di output quando non si desidera aumentare la velocità della linea. È possibile impedire la perdita di output solo se la perdita di output è una conseguenza di brevi picchi di dati. Se le perdite di output sono causate da un flusso costante ad alta velocità, non è possibile evitarle. Tuttavia, è possibile controllarli.

Tipi di accodamento coinvolti nella pianificazione dell'output sugli switch Catalyst 6500/6000

Caduta coda

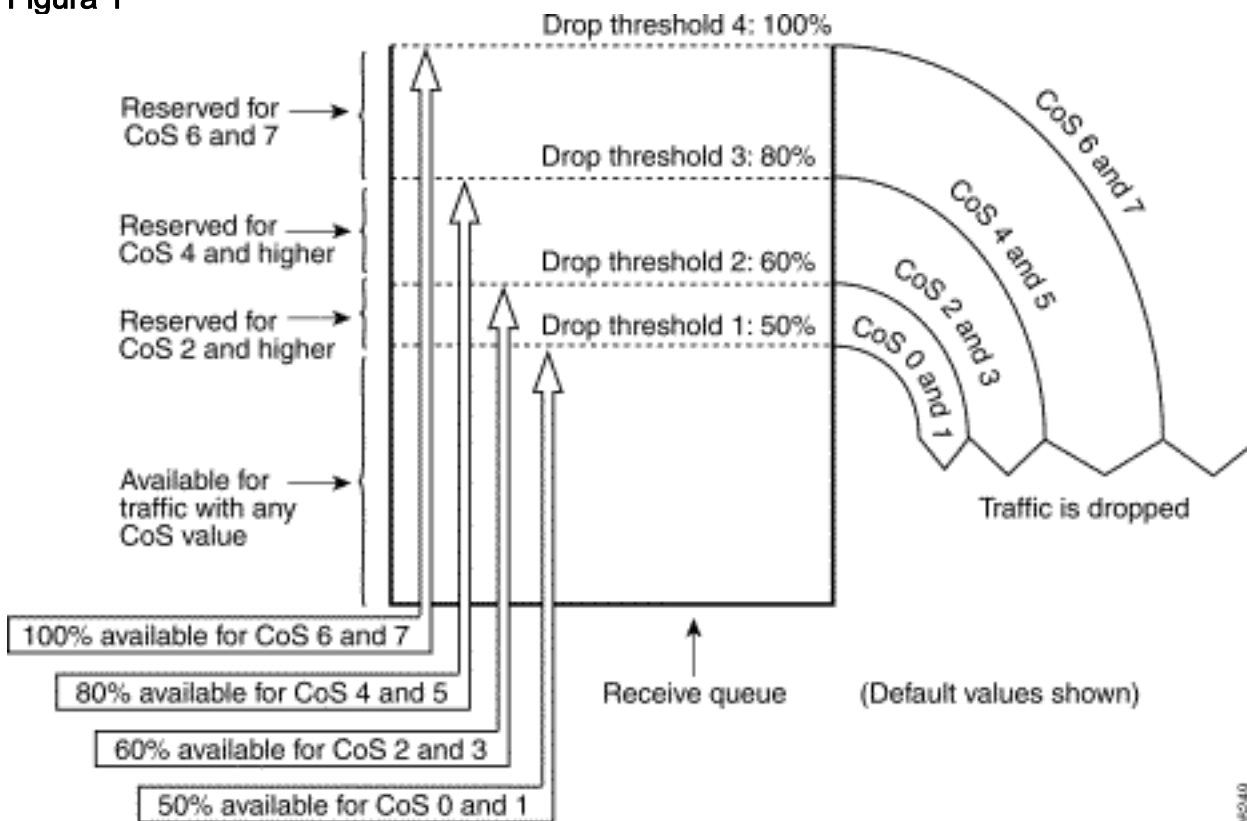
La caduta della coda è un meccanismo base per evitare le congestioni. Il servizio Tail Drop gestisce tutto il traffico in modo equo e non fa distinzioni tra le classi di servizio (CoS) quando le code iniziano a riempirsi durante i periodi di congestione. Quando la coda di output è piena ed è attivo il drop, i pacchetti vengono scartati finché non viene eliminata la congestione e la coda non è più piena. La perdita di velocità è il tipo più semplice di prevenzione delle congestioni e non tiene conto di alcun parametro QoS.

Catalyst 6000 ha implementato una versione avanzata di tail drop congestion che elimina tutti i pacchetti con un certo CoS quando si raggiunge una certa percentuale di riempimento del buffer. Con questa opzione è possibile definire una serie di soglie e associare un CoS a ciascuna soglia. Nell'esempio di questa sezione sono disponibili quattro possibili soglie. Le definizioni di ciascuna soglia sono le seguenti:

- La soglia 1 viene raggiunta quando viene riempito il 50% del buffer. CoS 0 e 1 sono assegnati a questa soglia.
- La soglia 2 viene raggiunta quando viene riempito il 60% del buffer. CoS 2 e 3 sono assegnati a questa soglia.
- La soglia 3 viene raggiunta quando viene riempito l'80% del buffer. A questa soglia sono assegnati i CoS 4 e 5.
- La soglia 4 viene raggiunta quando viene riempito il 100% del buffer. CoS 6 e 7 sono assegnati a questa soglia.

Nel diagramma della [Figura 1](#), tutti i pacchetti con un CoS di 0 o 1 vengono scartati se il buffer è pieno al 50%. Tutti i pacchetti con un valore CoS pari a 0, 1, 2 o 3 vengono scartati se i buffer sono riempiti al 60%. I pacchetti con un valore CoS di 6 o 7 vengono scartati quando i buffer sono completamente riempiti.

Figura 1



Nota: non appena il livello di riempimento del buffer scende al di sotto di una determinata soglia, i pacchetti con il CoS associato non vengono più scartati.

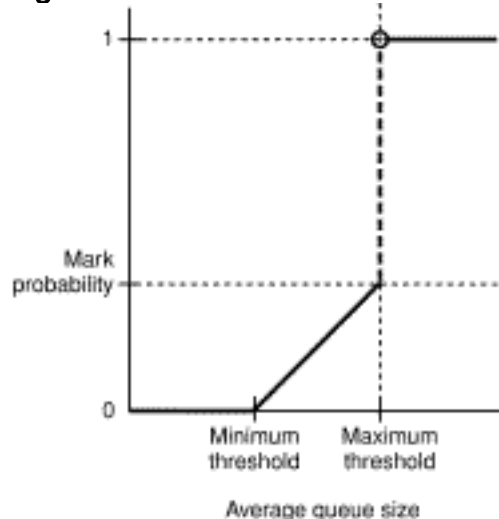
Random Early Detection e Weighted Random Early Detection

Il Weighted Random Early Detection (WRED) è un meccanismo di prevenzione della congestione che scarta casualmente i pacchetti con una certa precedenza IP quando i buffer raggiungono una determinata soglia di riempimento. WRED è una combinazione di queste due funzioni:

- Caduta
- RED (Random Early Detection)

RED non riconosce la precedenza o non riconosce CoS. RED utilizza una delle soglie singole quando il valore di soglia per il buffer è pieno. Il rosso inizia a rilasciare i pacchetti in modo casuale (ma non tutti, come nel caso del "tail drop") finché non viene raggiunta la soglia massima (max). Una volta raggiunta la soglia massima, tutti i pacchetti vengono scartati. La probabilità che un pacchetto venga scartato aumenta in modo lineare con l'aumento del riempimento del buffer al di sopra della soglia. Lo schema della [Figura 2](#) mostra la probabilità di perdita del pacchetto:

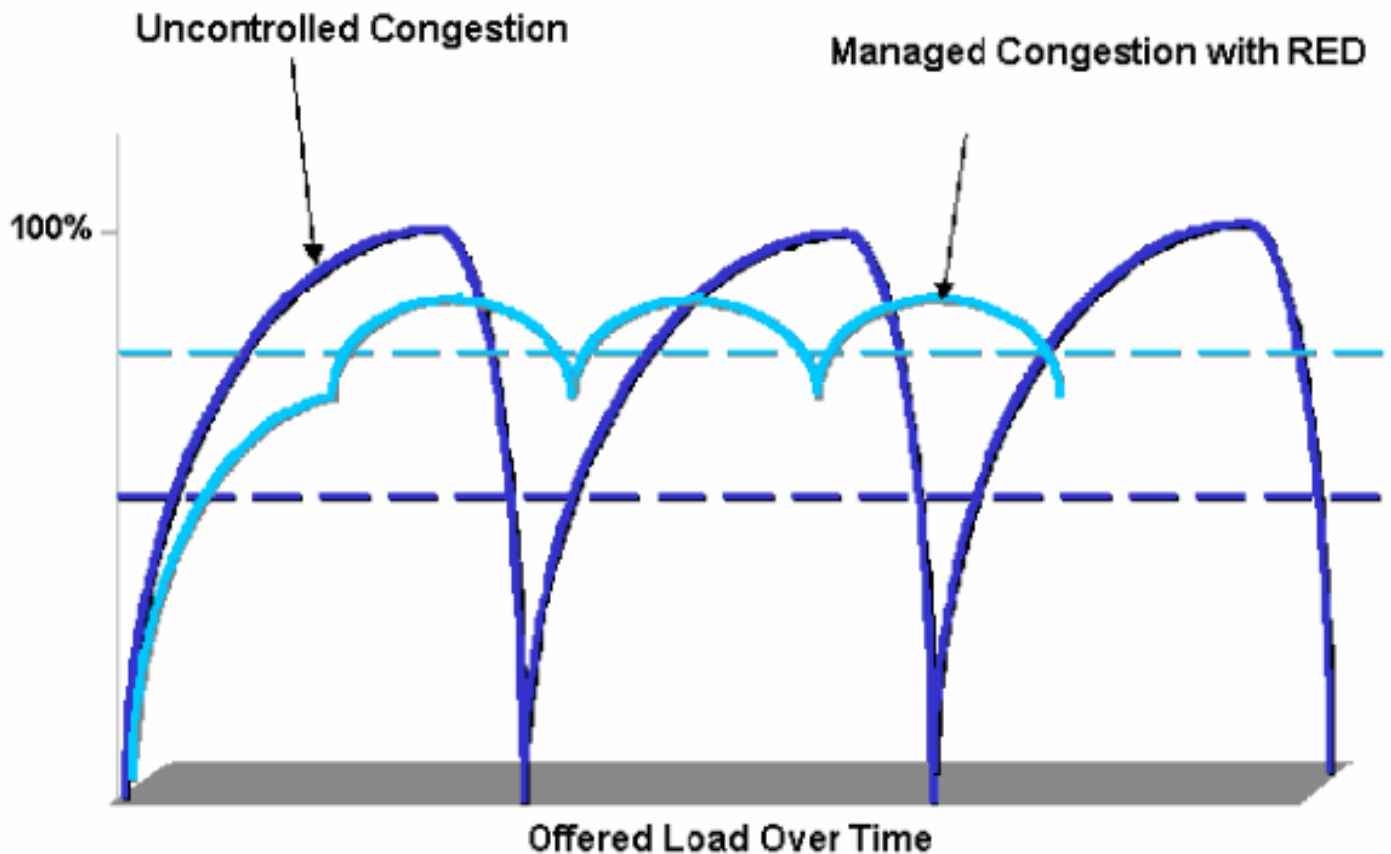
Figura 2 - Probabilità di eliminazione dei pacchetti



Nota: la probabilità di contrassegno in questo diagramma è regolabile in ROSSO, il che significa che la pendenza della probabilità di caduta lineare è regolabile.

RED e WRED sono meccanismi molto utili per prevenire la congestione del traffico basato su TCP. Per altri tipi di traffico, RED non è molto efficiente. Infatti, il protocollo RED sfrutta il meccanismo a finestre utilizzato dal protocollo TCP per gestire le congestioni. RED evita la tipica congestione che si verifica su un router quando più sessioni TCP attraversano la stessa porta del router. Il meccanismo è denominato sincronizzazione globale della rete. Lo schema della [Figura 3](#) mostra come il ROSSO abbia un effetto di livellamento sul carico:

Figura 3 - ROSSO per la prevenzione delle congestioni



Per ulteriori informazioni su come RED possa ridurre la congestione e regolare il traffico sul router, fare riferimento alla sezione [Come il router interagisce con il TCP](#) nel documento [Congestion Avoidance Overview](#).

WRED è simile a RED in quanto entrambi definiscono alcune soglie minime (min) e, quando queste soglie minime vengono raggiunte, i pacchetti vengono eliminati in modo casuale. WRED definisce inoltre determinate soglie massime e, quando tali soglie vengono raggiunte, tutti i pacchetti vengono eliminati. WRED è anche in grado di riconoscere il CoS, il che significa che uno o più valori del CoS vengono aggiunti a ciascuna coppia di soglia minima/soglia massima. Quando viene superata la soglia minima, i pacchetti vengono eliminati casualmente insieme al CoS assegnato. Considerare questo esempio con due soglie nella coda:

- CoS 0 e 1 vengono assegnati alla soglia minima 1 e alla soglia massima 1. La soglia minima 1 viene impostata sul 50% del riempimento del buffer e la soglia massima 1 viene impostata sull'80%.
- CoS 2 e 3 sono assegnati alla soglia minima 2 e alla soglia massima 2. La soglia minima 2 è impostata sul 70% del riempimento del buffer e la soglia massima 2 è impostata sul 100%.

Non appena il buffer supera la soglia minima 1 (50%), i pacchetti con CoS 0 e 1 iniziano a essere scartati in modo casuale. All'aumento dell'utilizzo del buffer, vengono scartati altri pacchetti. Se si raggiunge la soglia minima 2 (70%), i pacchetti con CoS 2 e 3 iniziano ad essere scartati casualmente.

Nota: in questa fase, la probabilità di perdita per i pacchetti con CoS 0 e 1 è molto più alta della probabilità di perdita per i pacchetti con CoS 2 o CoS 3.

Quando si raggiunge la soglia massima 2, i pacchetti con CoS 0 e 1 vengono scartati, mentre i pacchetti con CoS 2 e 3 continuano a essere scartati casualmente. Infine, quando si raggiunge il 100% (soglia massima 2), tutti i pacchetti con CoS 2 e 3 vengono scartati.

I diagrammi della [Figura 4](#) e della [Figura 5](#) illustrano un esempio di queste soglie:

Figura 4 - WRED con due set di soglie minime e massime (due servizi)

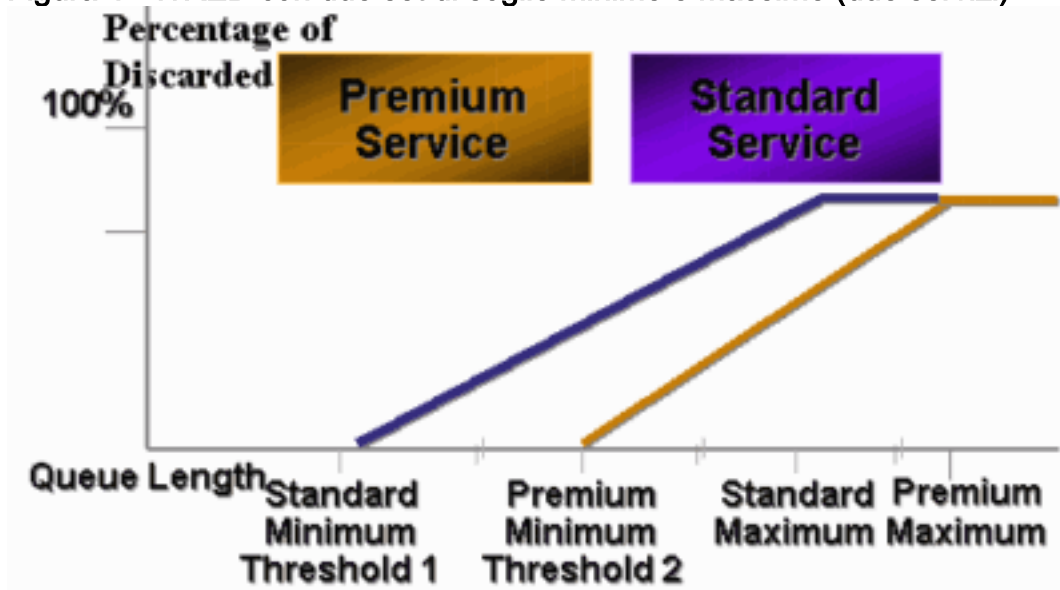
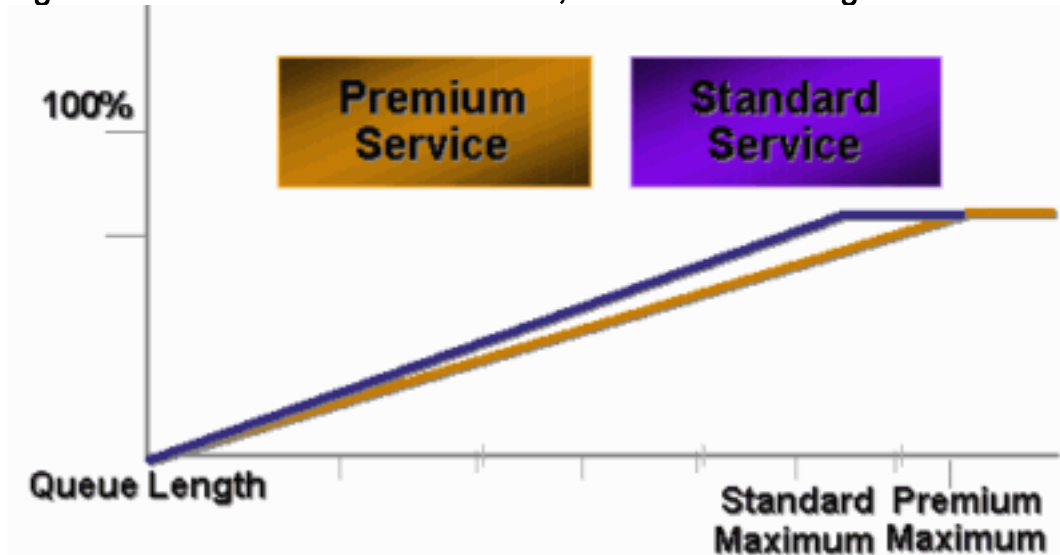


Figura 5 - WRED con due set di servizi, ma entrambe le soglie minime sono uguali a 0



La prima implementazione CatOS di WRED ha impostato solo la soglia massima, mentre la soglia minima è stata hardcoded allo 0%. La parte inferiore del diagramma nella [Figura 5](#) evidenzia il comportamento risultante.

Nota: la probabilità di perdita per un pacchetto è sempre diversa da null perché questa probabilità è sempre superiore alla soglia minima. Questo comportamento è stato corretto nel software versione 6.2 e successive.

[Round Robin ponderato](#)

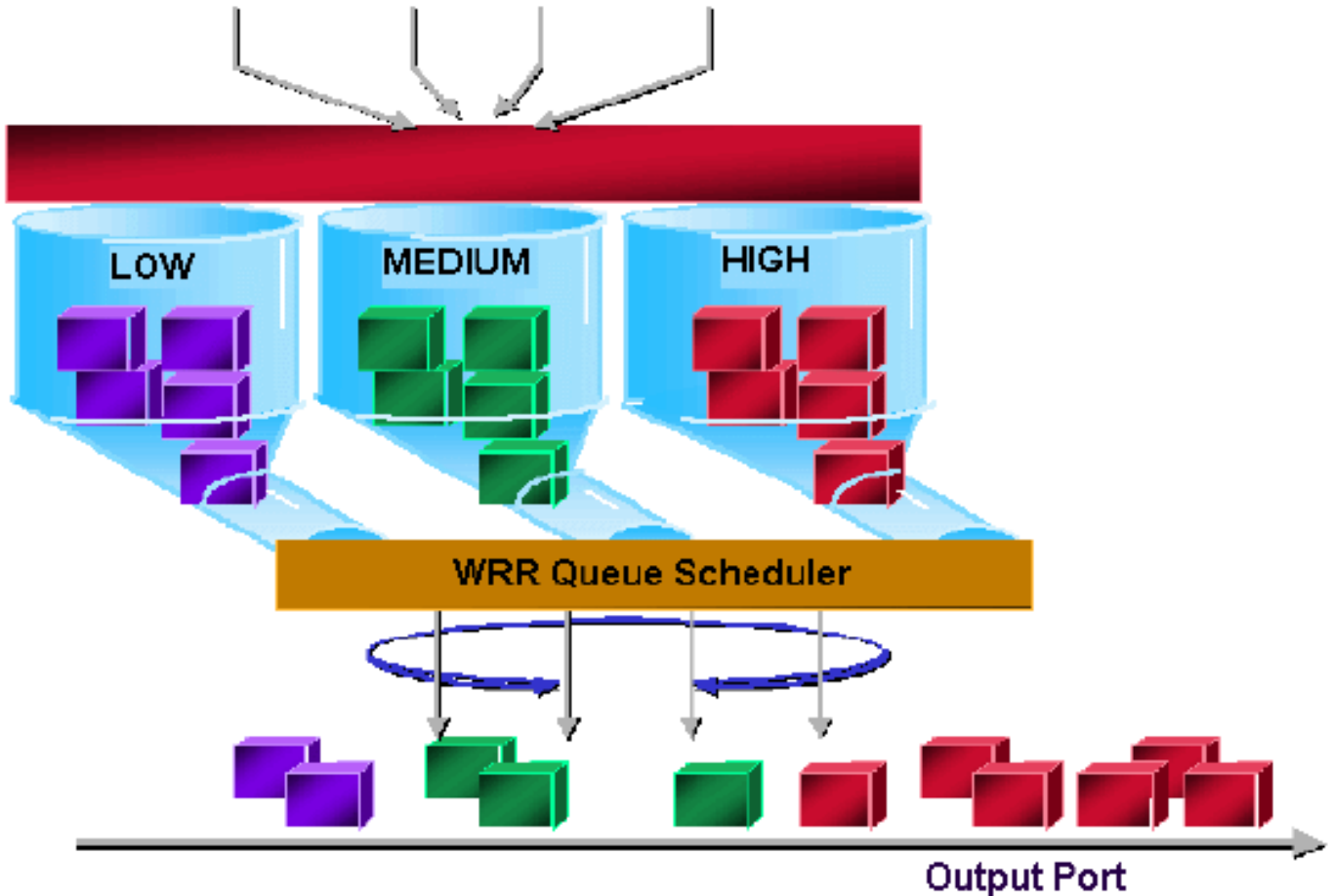
Il WRR (Weighted Round-Robin) è un altro meccanismo di programmazione dell'output su Catalyst 6000. WRR funziona tra due o più code. Le code per WRR vengono svuotate in modo round robin ed è possibile configurare il peso per ciascuna coda. Per impostazione predefinita, le porte hanno due code WRR su Catalyst 6000. Il valore predefinito è:

- Per servire la coda WRR con priorità alta il 70% del tempo

- Per servire la coda WRR con priorità bassa il 30% del tempo

Il diagramma della [Figura 6](#) mostra un WRR con tre code servite in modalità WRR. La coda ad alta priorità (pacchetti rossi) invia più pacchetti delle altre due code:

Figura 6 - Programmazione dell'uscita: WRR



Nota: la maggior parte delle schede di linea 6500 implementa il WRR per larghezza di banda. Questa implementazione di WRR per larghezza di banda significa che ogni volta che lo scheduler consente a una coda di trasmettere pacchetti, è consentito trasmettere un certo numero di byte. Questo numero di byte può rappresentare più di un pacchetto. Ad esempio, se si inviano 5120 byte di seguito, è possibile inviare tre pacchetti da 1518 byte, per un totale di 4554 byte. I byte in eccesso vengono persi ($5120 - 4554 = 566$ byte). Pertanto, con un peso eccessivo (come 1% per la coda 1 e 99% per la coda 2), il peso configurato esatto potrebbe non essere raggiunto. Questo mancato raggiungimento del peso esatto è spesso il caso dei pacchetti più grandi.

Alcune schede di linea di nuova generazione, come la 6548-RJ-45, superano questo limite attraverso l'implementazione del DWRR (deficit weighted round-robin). DWRR trasmette dalle code ma non priva la coda a bassa priorità. DWRR tiene traccia della coda a bassa priorità che si trova in fase di trasmissione e compensa nel ciclo successivo.

Coda con priorità rigorosa

Un altro tipo di coda in Catalyst 6000, una coda con priorità assoluta, viene sempre svuotata per prima. Non appena un pacchetto è presente nella coda di priorità rigida, il pacchetto viene inviato.

Le code WRR o WRED vengono controllate solo dopo lo svuotamento della coda con priorità assoluta. Dopo aver trasmesso ciascun pacchetto dalla coda WRR o WRED, la coda con priorità assoluta viene controllata e, se necessario, svuotata.

Nota: tutte le schede di linea con tipo di accodamento simile a 1p2q1t, 1p3q8t e 1p7q8t utilizzano DWRR. Le altre schede di linea utilizzano lo standard WRR.

Capacità di accodamento dell'output di diverse schede di linea su Catalyst 6000

Funzionalità comando show port

Se non si è certi della funzionalità di accodamento di una porta, è possibile usare il comando **show port capabilities**. Di seguito viene riportato l'output del comando su una scheda di linea WS-X6408-GBIC:

```
Model                WS-X6408-GBIC
Port                 4/1
Type                 No GBIC
Speed                1000
Duplex               full
Trunk encap type     802.1Q,ISL
Trunk mode           on,off,desirable,auto,nonegotiate
Channe               yes
Broadcast suppression percentage(0-100)
Flow control         receive-(off,on,desired),send-(off,on,desired)
Security             yes
MembershIP           static,dynamic
Fast start           yes
QOS scheduling       rx-(1q4t),tx-(2q2t)
CoS rewrite          yes
ToS rewrite          DSCP
UDLD                 yes
SPAN                 source,destination
COPS port group      none
```

Questa porta dispone di un tipo di output di accodamento denominato 2q2t.

Comprendere la funzionalità di accodamento di una porta

Sugli switch Catalyst 6500/6000 sono disponibili diversi tipi di code. Le tabelle in questa sezione potrebbero diventare incomplete quando vengono rilasciate nuove schede di linea. Le nuove schede di linea possono introdurre nuove combinazioni di accodamento. Per una descrizione aggiornata di tutte le code disponibili per i moduli degli switch Catalyst 6500/6000, fare riferimento alla sezione *Configurazione di QoS* per la versione CatOS della [documentazione del software Catalyst serie 6500](#).

Nota: Cisco Communication Media Module (CMM) non supporta tutte le funzionalità QoS. Per determinare le funzionalità supportate, consultare le note sulla versione del software in uso.

Nella tabella seguente viene illustrata la notazione dell'architettura QoS della porta:

Tx ¹ /R x ² ide	Notazion e coda	N. di code	Coda di priorità	N. di code WRR	N. e tipo di soglia per le code WRR
Tx	2q2t	2	—	2	2 drop di coda configurabile

Tx	1p2q2t	3	1	2	2 WRED configurabili
Tx	1p3q1t	4	1	3	1 WRED configurabile
Tx	1p2q1t	3	1	2	1 WRED configurabile
Rx	1q4t	1	—	1	4 drop di coda configurabile
Rx	1p1q4t	2	1	1	4 drop di coda configurabile
Rx	1p1q0t	2	1	1	Non configurabile
Rx	1p1q8t	2	1	1	8 WRED configurabili

¹ Tx = trasmissione.

² Rx = ricezione.

In questa tabella vengono elencati tutti i moduli e i tipi di coda sul lato Rx e Tx dell'interfaccia o della porta:

Modulo	Code Rx	Code Tx
WS-X6K-S2-PFC2	1p1q4t	1p2q2t
WS-X6K-SUP1A-2GE	1p1q4t	1p2q2t
WS-X6K-SUP1-2GE	1q4t	2q2t
WS-X6501-10GEX4	1p1q8t	1p2q1t
WS-X6502-10GE	1p1q8t	1p2q1t
WS-X6516-GBIC	1p1q4t	1p2q2t
WS-X6516-GE-TX	1p1q4t	1p2q2t
WS-X6416-GBIC	1p1q4t	1p2q2t
WS-X6416-GE-MT	1p1q4t	1p2q2t
WS-X6316-GE-TX	1p1q4t	1p2q2t
WS-X6408A-GBIC	1p1q4t	1p2q2t
WS-X6408-GBIC	1q4t	2q2t
WS-X6524-100FX-MM	1p1q0t	1p3q1t
WS-X6324-100FX-SM	1q4t	2q2t
WS-X6324-100FX-MM	1q4t	2q2t
WS-X624-100FX-MT	1q4t	2q2t
WS-X6548-RJ-21	1p1q0t	1p3q1t
WS-X6548-RJ-45	1p1q0t	1p3q1t
WS-X6348-RJ-21	1q4t	2q2t
WS-X6348-RJ21V	1q4t	2q2t
WS-X6348-RJ-45	1q4t	2q2t
WS-X6348-RJ-45V	1q4t	2q2t

WS-X6148-RJ-45V	1q4t	2q2t
WS-X6148-RJ21V	1q4t	2q2t
WS-X6248-RJ-45	1q4t	2q2t
WS-X6248A-TEL	1q4t	2q2t
WS-X6248-TEL	1q4t	2q2t
WS-X6024-10FL-MT	1q4t	2q2t

[Creazione di QoS sugli switch Catalyst 6500/6000](#)

Per produrre la funzionalità QoS, vengono utilizzati tre campi dello switch Catalyst 6500/6000:

- Precedenza IP: i primi tre bit del campo Type of service (ToS) nell'intestazione IP
- Il punto di codice dei servizi differenziati (DSCP): i primi sei bit del campo ToS nell'intestazione IP
- CoS - I tre bit utilizzati al livello 2 (L2) Questi tre bit fanno parte dell'intestazione ISL (Inter-Switch Link) o sono inclusi nel tag IEEE 802.1Q (dot1q). Non è presente alcun CoS all'interno di un pacchetto Ethernet senza tag.

[Meccanismo di programmazione dell'output su Catalyst 6500/6000](#)

Quando si invia un frame dal bus di dati da trasmettere, il CoS del pacchetto è l'unico parametro preso in considerazione. Il pacchetto passa quindi attraverso uno scheduler, che sceglie la coda in cui è inserito. Pertanto, tenere presente che la pianificazione dell'output e tutti i meccanismi illustrati in questo documento sono compatibili solo con CoS.

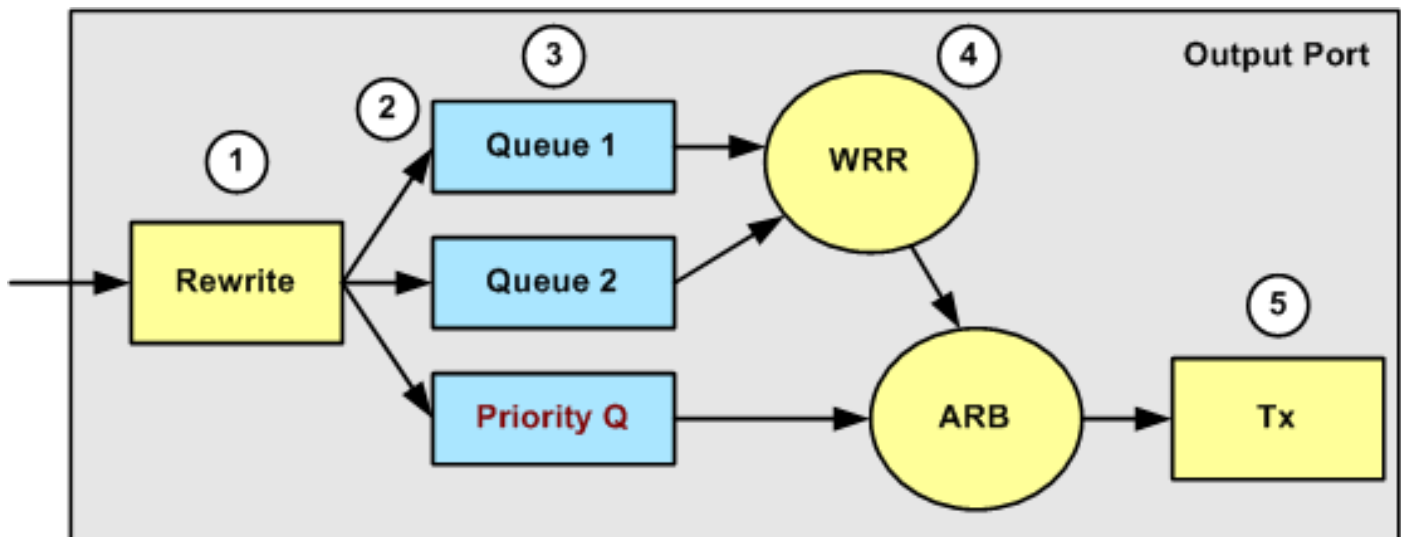
Per classificare il pacchetto, lo switch Catalyst 6500/6000 con un modulo Multilayer Switch Feature Card (MSFC) usa un DSCP interno. Gli switch Catalyst 6500/6000 configurati con QoS abilitato assegnano un valore DSCP quando la decisione di inoltrare viene presa a livello PFC. Questo DSCP viene assegnato a qualsiasi pacchetto, inclusi i pacchetti non IP, e viene mappato sul CoS per abilitare la pianificazione dell'output. È possibile configurare il mapping dai valori DSCP ai valori CoS sugli switch Catalyst 6500/6000. Se si lascia il valore predefinito, è possibile derivare il CoS dal DSCP. La formula è:

$$DSCP_value / 8$$

Inoltre, il valore DSCP viene mappato nel CoS del pacchetto in uscita, se il pacchetto è un pacchetto IP con tag ISL o dot1q (VLAN non nativa). Il valore DSCP viene scritto anche nel campo ToS dell'intestazione IP.

Lo schema della [Figura 7](#) mostra una coda 1p2q2t. Le code WRR vengono svuotate con l'utilità di pianificazione WRR. Inoltre, è presente un arbitro che controlla tra i pacchetti provenienti dalle code WRR per stabilire se contengono qualcosa nella coda con priorità rigorosa.

Figura 7



1. Il campo ToS (AS) viene riscritto nell'intestazione IP e il campo 802.1p/ISL CoS (CoS).
2. La coda di pianificazione e la soglia vengono selezionate in base al CoS, tramite una mappa configurabile.
3. Ogni coda ha dimensioni e soglie configurabili e alcune code hanno WRED.
4. La rimozione dalla coda utilizza WRR tra due code.
5. L'incapsulamento in uscita può essere dot1q, ISL o none.

[Pianificazione della configurazione, del monitoraggio e dell'output sugli switch Catalyst 6500/6000](#)

[Configurazione predefinita per QoS sugli switch Catalyst 6500/6000](#)

In questa sezione viene fornito un output di esempio dalla configurazione QoS predefinita su uno switch Catalyst 6500/6000, oltre a informazioni sul significato di questi valori e sulla modalità di tuning dei valori.

QoS è disabilitato per impostazione predefinita quando si usa questo comando:

```
set qos disable
```

I comandi in questo elenco mostrano l'assegnazione predefinita per ogni CoS in una porta 2q2t. Alla coda 1 è assegnato il CoS 0 e 1 alla prima soglia e il CoS 2 e il CoS 3 alla seconda soglia. Alla coda 2 sono assegnati i livelli CoS 4 e 5 alla prima soglia e i livelli CoS 6 e 7 alla seconda soglia:

```
set qos map 2q2t tx 1 1 cos 0
```

```
set qos map 2q2t tx 1 1 cos 1
```

```
set qos map 2q2t tx 1 2 cos 2
```

```
set qos map 2q2t tx 1 2 cos 3
```

```
set qos map 2q2t tx 2 1 cos 4
```

```
set qos map 2q2t tx 2 1 cos 5
```

```
set qos map 2q2t tx 2 2 cos 6
```

```
set qos map 2q2t tx 2 2 cos 7
```

Questi comandi visualizzano il livello di soglia per impostazione predefinita su una porta 2q2t per ciascuna coda:

```
set qos drop-threshold 2q2t tx queue 1 80 100
```

```
set qos drop-threshold 2q2t tx queue 2 80 100
```

È possibile assegnare il peso predefinito a ciascuna delle code WRR. Utilizzare questo comando per assegnare i pesi predefiniti per la coda 1 e la coda 2:

Nota: la coda con priorità bassa viene servita il 5/260% del tempo, mentre la coda con priorità alta viene servita il 25/260% del tempo.

```
set qos wrr 2q2t 5 255
```

La disponibilità totale del buffer viene suddivisa tra le due code. La coda a bassa priorità è correttamente assegnata all'80% dei buffer disponibili perché è la coda per la quale è più probabile che i pacchetti vengano memorizzati nel buffer e rimangano in attesa per un certo periodo di tempo. Per definire la disponibilità, usare questo comando:

```
set qos txq-ratio 2q2t 80 20
```

È possibile visualizzare impostazioni simili per la porta 1p2q2t in questa configurazione:

```
set qos map 1p2q2t tx 1 1 cos 0
```

```
set qos map 1p2q2t tx 1 1 cos 1
```

```
set qos map 1p2q2t tx 1 2 cos 2
```

```
set qos map 1p2q2t tx 1 2 cos 3
```

```
set qos map 1p2q2t tx 2 1 cos 4
```

```
set qos map 1p2q2t tx 3 1 cos 5
```

```
set qos map 1p2q2t tx 2 1 cos 6
```

```
set qos map 1p2q2t tx 2 2 cos 7
```

```
set qos wrr 1p2q2t 5 255
```

```
set qos txq-ratio 1p2q2t 70 15 15
```

```
set qos wred 1p2q2t tx queue 1 80 100
```

```
set qos wred 1p2q2t tx queue 2 80 100
```

Nota: per impostazione predefinita, CoS 5 (traffico voce) viene assegnato alla coda con priorità rigida.

Configurazione

Il primo passaggio della configurazione è abilitare QoS. Tenere presente che QoS è disabilitato per impostazione predefinita. Quando QoS è disabilitato, il mapping CoS è irrilevante. C'è una sola coda che viene servita come FIFO, e tutti i pacchetti vengono scartati lì.

```
bratan> (enable) set qos enable
```

```
QoS is enabled
```

```
bratan> (enable) show qos status
```

```
QoS is enabled on this switch
```

Il valore CoS deve essere assegnato alla coda o alla soglia per tutti i tipi di coda. Il mapping definito per un tipo di porta 2q2t non viene applicato ad alcuna porta 1p2q2t. Inoltre, il mapping creato per 2q2t viene applicato a tutte le porte con un meccanismo di coda 2q2t. Immettere questo comando

```
set qos map queue_type tx Q_number threshold_number cos value
```

Nota: le code sono sempre numerate in modo da iniziare con la coda con la priorità più bassa possibile e terminare con la coda con priorità assoluta disponibile. Di seguito è riportato un esempio:

- La coda 1 è la coda WRR a bassa priorità
- La coda 2 è la coda WRR ad alta priorità
- La coda 3 è la coda con priorità assoluta

È necessario ripetere questa operazione per tutti i tipi di code. In caso contrario, viene mantenuta l'assegnazione CoS predefinita. Di seguito è riportato un esempio per 1p2q2t:

Configurazione

```
set qos map 1p2q2t tx 1 1 cos 0
```

```
!--- This is the low-priority WRR queue threshold 1, CoS 0 and 1. set qos map 1p2q2t tx 1 1 cos 1 and 1
```

```
set qos map 1p2q2t tx 1 2 cos 2
```

```
!--- This is the low-priority WRR queue threshold 2, CoS 2 and 3. set qos map 1p2q2t tx 1 2 cos 3 and 3
```

```
set qos map 1p2q2t tx 2 1 cos 4
```

```
!--- This is the high-priority WRR queue threshold 1, CoS 4. set qos map 1p2q2t tx 3 1 cos 5
!--- This is the strict priority queue, CoS 5. set qos map 1p2q2t tx 2 1 cos 6
!--- This is the high-priority WRR queue threshold 2, CoS 6. set qos map 1p2q2t tx 2 2 cos 7 and
7
```

Output della console

```
tamer (enable) set qos map 1p2q2t tx 1 1 cos 0
```

QoS tx priority queue and threshold mapped to cos successfully

È necessario configurare il peso WRR per le due code WRR. Immettere questo comando

```
set qos wrr Q_type weight_1 weight_2
```

Weight_1 si riferisce alla coda 1, che deve essere la coda WRR a bassa priorità. *Weight_1* deve sempre essere inferiore a *weight_2*. Il peso può assumere qualsiasi valore compreso tra 1 e 255. È possibile assegnare la percentuale con le seguenti formule:

- Coda 1:

$$\text{weight}_1 / (\text{weight}_1 + \text{weight}_2)$$

- Coda 2:

$$\text{weight}_2 / (\text{weight}_1 + \text{weight}_2)$$

È inoltre necessario definire il peso per i vari tipi di code. Il peso non deve essere necessariamente lo stesso. Ad esempio, per 2q2t, dove la coda 1 viene servita il 30% del tempo e la coda 2 viene servita il 70% del tempo, è possibile eseguire questo comando per definire il peso:

```
set qos wrr 2q2t 30 70
```

```
!--- This ensures that the high-priority WRR queue is served 70 percent of the time !--- and
that the low-priority WRR queue is served 30 percent of the time.
```

Output della console

```
tamer (enable) set qos wrr 2q2t 30 70
```

QoS wrr ratio is set successfully

È inoltre necessario definire il rapporto della coda di trasmissione, che si riferisce al modo in cui i buffer vengono suddivisi tra le diverse code. Immettere questo comando

```
set qos txq-ratio port_type queue1_val queue2_val ... queueN_val
```

Nota: se si dispone di tre code (1p2q2t), è necessario impostare la coda WRR con priorità alta e la coda con priorità assoluta allo stesso livello per motivi hardware.

Configurazione

```
set qos txq-ratio 1p2q2t 70 15 15
```

!--- This gives 70 percent of the buffer of all 1p2q2t ports to the low-priority WRR !--- queue and gives 15 percent to each of the other two queues. set qos txq-ratio 2q2t 80 20
!--- This gives 80 percent of the buffer to the low-priority queue, !--- and gives 20 percent of the buffer to the high-priority queue.

Output della console

```
tamer (enable) set qos txq-ratio 1p2q2t 70 15 20
```

Queue ratio values must be in range of 1-99 and add up to 100
Example: set qos txq-ratio 2q2t 20 80

```
tamer (enable) set qos txq-ratio 1p2q2t 70 30 30
```

Queue ratio values must be in range of 1-99 and add up to 100
Example: set qos txq-ratio 2q2t 20 80

```
tamer (enable) set qos txq-ratio 1p2q2t 80 10 10
```

QoS txq-ratio is set successfully

Come illustrato nell'output di questa console, la somma dei valori della coda deve essere 100. Lasciare la parte più grande dei buffer per la coda WRR a bassa priorità perché questa coda richiede il maggior numero di buffer. Le altre code vengono servite con priorità più alta.

L'ultimo passaggio consiste nel configurare il livello di soglia per la coda WRED o per la coda di rilascio. Utilizzare i seguenti comandi:

```
set qos wred port_type [tx] queue q_num thr1 thr2 ... thrn
```

```
set qos drop-threshold port_type tx queue q_num thr1 ... thr2
```

Configurazione

```
set qos drop-threshold 2q2t tx queue 1 50 80
```

!--- For low-priority queues in the 2q2t port, the first threshold is defined at 50 !--- percent and the second threshold is defined at 80 percent of buffer filling. set qos drop-threshold 2q2t tx queue 2 40 80

!--- For high-priority queues in the 2q2t port, the first threshold is defined at 40 !--- percent and the second threshold is defined at 80 percent of buffer filling. set qos wred 1p2q2t tx queue 1 50 90

!--- The commands for the 1p2q2t port are identical. set qos wred 1p2q2t tx queue 2 40 80

Output della console

```
tamer (enable) set qos drop-threshold 2q2t tx queue 1 50 80
```

Transmit drop thresholds for queue 1 set at 50% 80%

```
tamer (enable) set qos drop-threshold 2q2t tx queue 2 40 80
```

Transmit drop thresholds for queue 2 set at 40% 80%

```
tamer (enable) set qos wred lp2q2t tx queue 1 50 90
```

WRED thresholds for queue 1 set to 50 and 90 on all WRED-capable lp2q2t ports

```
tamer (enable) set qos wred lp2q2t tx queue 2 40 80
```

WRED thresholds for queue 2 set to 40 and 80 on all WRED-capable lp2q2t ports

Il comando **set qos wred lp2q2t tx queue 2 40 80** funziona con il CoS per il mapping delle soglie. Ad esempio, quando si immettono i comandi nell'elenco seguente, si verifica che sulla porta 1p2q2t nella direzione di trasmissione i pacchetti con CoS 0, 1, 2 e 3 vengano inviati alla prima coda (la coda WRR inferiore). Quando i buffer della coda sono riempiti al 50%, WRED inizia a rilasciare i pacchetti con CoS 0 e 1. I pacchetti con CoS 2 e 3 vengono scartati solo quando i buffer della coda sono riempiti al 90%.

```
set qos map lp2q2t tx 1 1 cos 0
```

```
set qos map lp2q2t tx 1 1 cos 1
```

```
set qos map lp2q2t tx 1 2 cos 2
```

```
set qos map lp2q2t tx 1 2 cos 3
```

```
set qos wred lp2q2t tx queue 1 50 90
```

[Monitorare la pianificazione dell'output e verificare la configurazione](#)

Un semplice comando da utilizzare per verificare la configurazione di runtime corrente per la pianificazione dell'output di una porta è **show qos info runtime *mod/porta***. Il comando visualizza le seguenti informazioni:

- Tipo di accodamento sulla porta
- Mapping di CoS alle diverse code e soglie
- Condivisione del buffer
- Il peso WRR

In questo esempio, i valori sono WRR 20% per la coda 1 e WRR 80% per la coda 2:

```
tamer (enable) show qos info runtime 1/1
```

Run time setting of QoS:

QoS is enabled

Policy Source of port 1/1: Local

Tx port type of port 1/1 : lp2q2t

Rx port type of port 1/1 : lp1q4t

Interface type: port-based

ACL attached:

The qos trust type is set to untrusted

Default CoS = 0

Queue and Threshold Mapping for lp2q2t (tx):

Queue	Threshold	CoS
1	1	0 1
1	2	2 3
2	1	4 6


```

2          2          7
3          1          5

```

Queue and Threshold Mapping for 1p1q4t (rx):

All packets are mapped to a single queue

Rx drop thresholds:

Rx drop thresholds are disabled

Tx drop thresholds:

Tx drop-thresholds feature is not supported for this port type

Tx WRED thresholds:

Queue # Thresholds - percentage (* abs values)

```

-----
1          80% (249088 bytes) 100% (311168 bytes)
2          80% (52480 bytes) 100% (61440 bytes)

```

Queue Sizes:

Queue # Sizes - percentage (* abs values)

```

-----
1          70% (311296 bytes)
2          15% (65536 bytes)
3          15% (65536 bytes)

```

WRR Configuration of ports with speed 1000Mbps:

Queue # Ratios (* abs values)

```

-----
1          20 (5120 bytes)
2          80 (20480 bytes)

```

(*) Runtime information may differ from user configured setting due to hardware granularity.

tamer (enable)

Nell'esempio successivo, i pesi WRR non sono il valore predefinito 1. I pesi sono stati impostati sui valori 20 per la coda 1 e 80 per la coda 2. In questo esempio viene utilizzato un generatore di traffico per inviare 2 Gb di traffico a Catalyst 6000. Questi 2 Gb di traffico dovrebbero uscire dalla porta 1/1. Poiché la porta 1/1 ha una sottoscrizione eccessiva, molti pacchetti vengono scartati (1 Gbps). Il comando **show mac** mostra una notevole perdita di output:

tamer (enable) **show mac 1/1**

```

Port          Rcv-Unicast          Rcv-Multicast          Rcv-Broadcast
-----
1/1           0                    1239                   0

Port          Xmit-Unicast          Xmit-Multicast          Xmit-Broadcast
-----
1/1          73193601             421                    0

Port          Rcv-Octet            Xmit-Octet
-----
1/1          761993              100650803690

MAC          Dely-Exced          MTU-Exced          In-Discard          Out-Discard
-----
1/1          0                   -                   0                   120065264

```

Last-Time-Cleared

```

-----
Fri Jan 12 2001, 17:37:43

```

Prendiamo in considerazione i pacchetti che vengono scartati. Il modello di traffico suggerito viene suddiviso in questo modo:

- 1 Gb di traffico con IP precedence 0
- 250 MB di traffico con IP Precence 4

- 250 MB di traffico con IP Precedence 5
- 250 MB di traffico con precedenza IP 6
- 250 MB di traffico con IP precedence 7

Secondo la mappatura CoS, questo traffico viene inviato:

- 1 Gb di traffico verso la coda 1 soglia 1
- 0 MB di traffico verso la coda 1 soglia 2
- 500 MB di traffico verso la coda 2 soglia 1
- 250 MB di traffico verso la coda 2 soglia 2
- 250 MB di traffico verso la coda 3 (coda con priorità assoluta)

Lo switch deve considerare attendibile il traffico ricevuto in modo che la precedenza IP in ingresso venga mantenuta nello switch e venga utilizzata per il mapping al valore CoS per la pianificazione dell'output.

Nota: la precedenza IP predefinita per il mapping CoS è precedenza IP uguale a CoS.

Utilizzare il comando **show qos stat 1/1** per verificare i pacchetti ignorati e la percentuale approssimativa:

- A questo punto, nessun pacchetto viene scartato nella coda 3 (CoS 5).
- Il 91,85% dei pacchetti scartati sono pacchetti CoS 0 nella coda 1.
- L'8% dei pacchetti scartati sono CoS 4 e 6 nella coda 2, soglia 1.
- Lo 0,15% dei pacchetti scartati è CoS 7 nella coda 2, soglia 2.

Questo output illustra l'utilizzo del comando:

```
tamer (enable) show qos stat 1/1
```

```
Tx port type of port 1/1 : 1p2q2t
```

```
Q3T1 statistics are covered by Q2T2.
```

```
Q #      Threshold #:Packets dropped
-----
1        1:110249298 pkts, 2:0 pkts
2        1:9752805 pkts, 2:297134 pkts
3        1:0 pkts
```

```
Rx port type of port 1/1 : 1p1q4t
```

```
Rx drop threshold counters are disabled for untrusted ports
```

```
Q #      Threshold #:Packets dropped
-----
1        1:0 pkts, 2:0 pkts, 3:0 pkts, 4:0 pkts
2        1:0 pkts
```

Se si ripristina il valore predefinito del peso WRR dopo aver cancellato i contatori, solo l'1% dei pacchetti scartati si troverà nella coda 2, invece dell'8% visualizzato in precedenza:

Nota: il valore predefinito è 5 per la coda 1 e 255 per la coda 2.

```
tamer (enable) show qos stat 1/1
```

```
TX port type of port 1/1 : 1p2q2t
```

```
Q3T1 statistics are covered by Q2T2
```

```
Q #      Threshold #:Packets dropped
-----
1        1:2733942 pkts, 2:0 pkts
2        1:28890 pkts, 2:6503 pkts
```

```

3          1:0 pkts
Rx port type of port 1/1 : lp1q4t
Rx drop threshold counters are disabled for untrusted ports
Q #          Threshold #:Packets dropped
---          -----
1          1:0 pkts, 2:0 pkts, 3:0 pkts, 4:0 pkts
2          1:0 pkts

```

Utilizzare la programmazione dell'output per ridurre il ritardo e l'instabilità

L'esempio della sezione [Monitorare la pianificazione dell'output e verificare la configurazione](#) mostra i vantaggi dell'implementazione della pianificazione dell'output, che evita una perdita di traffico VoIP o mission-critical in caso di sottoscrizione eccessiva della porta di output.

L'oversubscription si verifica raramente in una rete normale, in particolare su un collegamento Gigabit. In genere, l'iscrizione in eccesso si verifica solo durante i periodi di picco del traffico o durante i picchi di traffico in un periodo di tempo molto breve.

Anche senza sottoscrizioni eccessive, la pianificazione dell'output può essere di grande aiuto in una rete in cui QoS è implementato end-to-end. La programmazione dell'output aiuta a ridurre il ritardo e l'instabilità. In questa sezione vengono forniti esempi di come la programmazione dell'output può contribuire a ridurre il ritardo e l'instabilità.

Riduci ritardo

Il ritardo di un pacchetto è incrementato dal tempo "perso" nel buffer di ciascuno switch durante l'attesa della trasmissione. Ad esempio, un piccolo pacchetto vocale con un CoS di 5 viene inviato da una porta durante un backup o un trasferimento di file di grandi dimensioni. Se non si dispone di QoS per la porta di output e si presume che il pacchetto voce piccolo venga inserito in coda dopo 10 pacchetti grandi da 1500 byte, è possibile calcolare facilmente il tempo di velocità Gigabit necessario per trasmettere i 10 pacchetti grandi:

`(10 × 1500 × 8) = 120,000 bits that are transmitted in 120 microseconds`

Se il pacchetto deve attraversare otto o nove switch mentre attraversa la rete, può verificarsi un ritardo di circa 1 ms. Questa quantità conta solo i ritardi nella coda di output dello switch che viene attraversato nella rete.

Nota: se è necessario accodare gli stessi 10 pacchetti di grandi dimensioni su un'interfaccia a 10 Mbps (ad esempio, con un telefono IP e un PC connesso), il ritardo introdotto è:

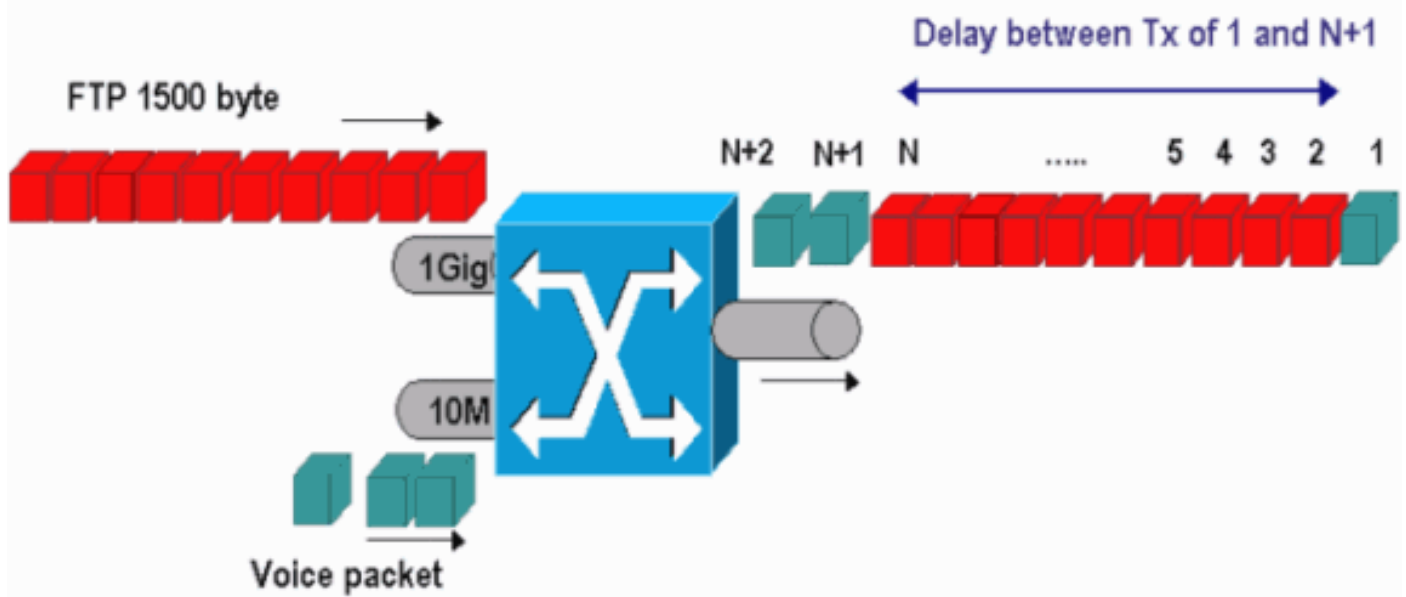
`(10 × 1500 × 8) = 120,000 bits that are transmitted in 12 ms`

L'implementazione della pianificazione dell'output assicura che i pacchetti voce con un CoS di 5 vengano inseriti nella coda con priorità rigorosa. Questo posizionamento assicura che questi pacchetti vengano inviati prima di qualsiasi pacchetto con un CoS inferiore a 5, riducendo i ritardi.

Riduci variazione

Un altro vantaggio importante dell'implementazione della programmazione dell'output è la riduzione del tremolio. Il jitter è la variazione di ritardo osservata per i pacchetti all'interno dello stesso flusso. Il diagramma della [Figura 8](#) mostra uno scenario di esempio in cui la programmazione dell'output può ridurre il jitter:

Figura 8



In questo scenario, una singola porta di output deve inviare due flussi:

- Un flusso vocale in ingresso su una porta Ethernet a 10 Mbps
- Un flusso FTP in ingresso su un uplink Ethernet a 1 Gbps

Entrambi i flussi lasciano lo switch attraverso la stessa porta di output. In questo esempio viene mostrato ciò che può accadere senza l'utilizzo della programmazione dell'output. Tutti i pacchetti di dati di grandi dimensioni possono essere interlacciati tra due pacchetti voce, creando un effetto jitter nella ricezione del pacchetto voce dallo stesso flusso. Il ritardo tra la ricezione del pacchetto n e il pacchetto $n+1$ è maggiore in quanto lo switch trasmette il pacchetto di dati di grandi dimensioni. Tuttavia, il ritardo tra $n+1$ e $n+2$ è trascurabile. Ciò determina uno jitter nel flusso del traffico vocale. È possibile evitare questo problema utilizzando una coda con priorità rigida. Verificare che il valore CoS dei pacchetti voce sia mappato alla coda con priorità rigida.

[Informazioni correlate](#)

- [Pianificazione dell'uscita QoS sugli switch Catalyst serie 6500/6000 con software Cisco IOS](#)
- [Qualità del servizio sugli switch Catalyst serie 6000](#)
- [Pagine di supporto dei prodotti LAN](#)
- [Pagina di supporto dello switching LAN](#)
- [Documentazione e supporto tecnico – Cisco Systems](#)