

Programação da saída de QoS nos Switches da série Catalyst 6500/6000 executando o Software do sistema CatOS

Contents

[Introduction](#)

[Prerequisites](#)

[Requirements](#)

[Componentes Utilizados](#)

[Conventions](#)

[Informações de Apoio](#)

[Quedas da fila da saída](#)

[Tipos de Enfileiramento Envolvidos na Programação de Saída no Catalyst 6500/6000](#)

[Queda traseira](#)

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

[Rodízio ponderado](#)

[Fila de prioridade estrita](#)

[Recurso de enfileiramento de saída de diferentes placas de linha no Catalyst 6000](#)

[Recursos do comando show port](#)

[Entender o recurso de enfileiramento de uma porta](#)

[Criar QoS no Catalyst 6500/6000](#)

[Mecanismo de programação de saída no Catalyst 6500/6000](#)

[Configuração, monitoramento e programação de saída no Catalyst 6500/6000](#)

[Configuração padrão para QoS no Catalyst 6500/6000](#)

[Configuração](#)

[Monitorar o Agendamento de Saída e Verificar a Configuração](#)

[Usar o Agendamento de Saída para Reduzir o Atraso e o Jitter](#)

[Reduza o atraso](#)

[Reduzir o atraso de sincronismo](#)

[Informações Relacionadas](#)

Introduction

O uso da programação de saída garante que o tráfego importante não caia no caso de um grande excesso de assinaturas. Este documento discute todas as técnicas e algoritmos envolvidos na programação de saída nos switches das séries Cisco Catalyst 6500/6000 que executam o software Catalyst OS (CatOS). Este documento também fornece uma breve visão geral da capacidade de enfileiramento dos switches Catalyst 6500/6000 e como configurar os diferentes parâmetros diferentes da programação de saída.

Observação: se você executar o Cisco IOS® Software em seus Catalyst 6500/6000, consulte [Programação de Saída de QoS nos Catalyst 6500/6000 Series Switches executando o Cisco IOS System Software](#) para obter mais informações.

Prerequisites

Requirements

Não existem requisitos específicos para este documento.

Componentes Utilizados

Os exemplos neste documento foram criados a partir de um Catalyst 6000 com um Supervisor Engine 1A e um Policy Feature Card (PFC). Mas os exemplos também são válidos para um Supervisor Engine 2 com um PFC2 ou para um Supervisor Engine 720 com um PFC3.

The information in this document was created from the devices in a specific lab environment. All of the devices used in this document started with a cleared (default) configuration. If your network is live, make sure that you understand the potential impact of any command.

Conventions

Consulte as [Convenções de Dicas Técnicas da Cisco para obter mais informações sobre convenções de documentos](#).

Informações de Apoio

Quedas da fila da saída

As quedas de saída são causadas por uma interface congestionada. Uma causa comum disso pode ser o tráfego de um link de largura de banda alta que está sendo comutado para um link de largura de banda mais baixa ou o tráfego de vários links de entrada que estão sendo comutados para um único link de saída.

Por exemplo, se uma grande quantidade de tráfego intermitente entra em uma interface gigabit e é comutada para uma interface de 100 Mbps, isso pode fazer com que quedas de saída aumentem na interface de 100 Mbps. Isso ocorre porque a fila de saída nessa interface está sobrecarregada pelo excesso de tráfego devido à incompatibilidade de velocidade entre as larguras de banda de entrada e saída. A taxa de tráfego na interface de saída não pode aceitar todos os pacotes que devem ser enviados.

A solução decisiva para resolver o problema é aumentar a velocidade da linha. No entanto, há maneiras de evitar, diminuir ou controlar quedas de saída quando você não deseja aumentar a velocidade da linha. Você pode evitar quedas de saída somente se quedas de saída forem uma consequência de rajadas curtas de dados. Se as quedas de saída forem causadas por um fluxo constante de alta taxa, você não poderá impedir as quedas. Entretanto, você pode controlá-los.

Tipos de Enfileiramento Envolvidos na Programação de Saída no

Catalyst 6500/6000

Queda traseira

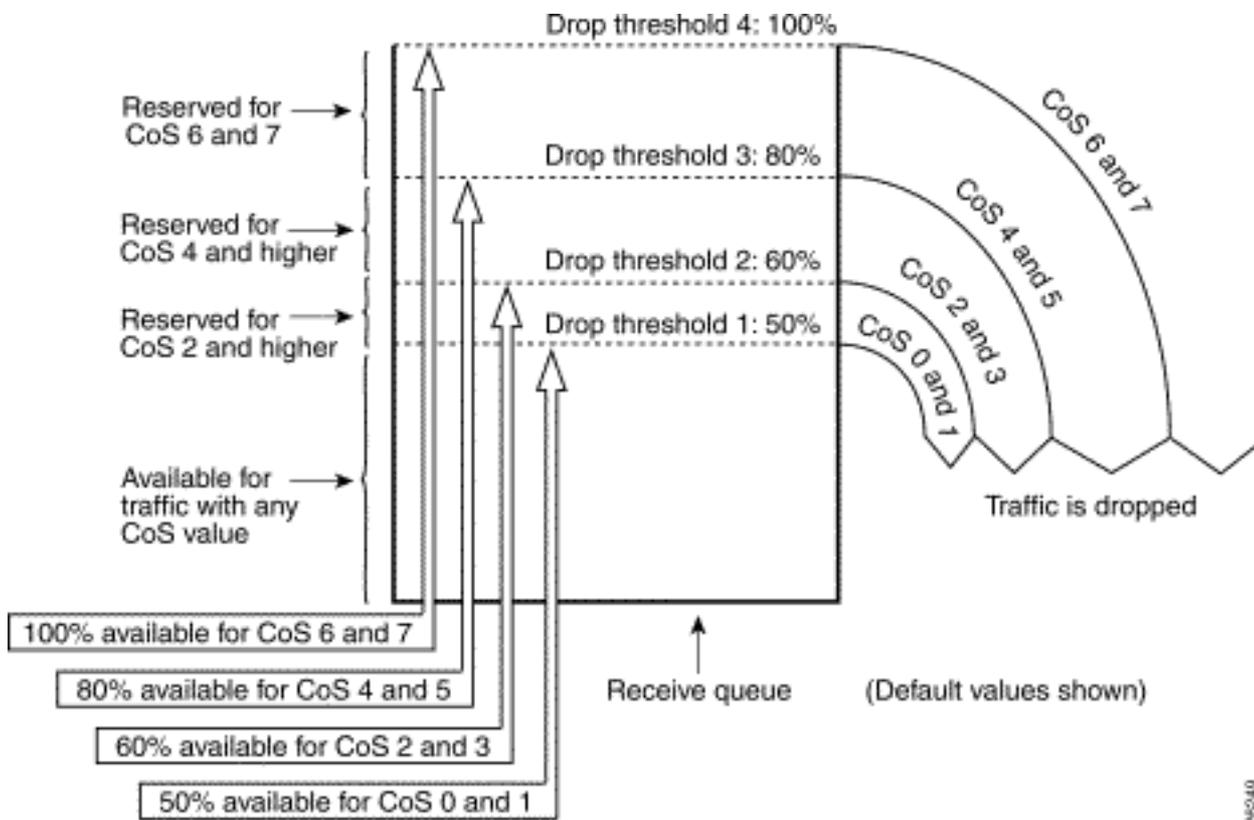
A queda traseira é um mecanismo básico de prevenção de congestionamento. A queda traseira trata todo o tráfego igualmente e não diferencia entre classes de serviço (CoS) quando as filas começam a ser preenchidas durante períodos de congestionamento. Quando a fila de saída está cheia e a queda traseira está em vigor, os pacotes são descartados até que o congestionamento seja eliminado e a fila não esteja mais cheia. A queda traseira é o tipo mais básico de prevenção de congestionamento e não leva em conta nenhum parâmetro de QoS.

O Catalyst 6000 implementou uma versão avançada de prevenção de congestionamento de queda traseira que descarta todos os pacotes com um determinado CoS quando um determinado percentual de preenchimento de buffer é alcançado. Com queda traseira ponderada, você pode definir um conjunto de limiares e associar um CoS a cada limite. No exemplo desta seção, há quatro limiares possíveis. As definições de cada limiar são:

- O limite 1 é atingido quando 50% do buffer é preenchido. Os CoS 0 e 1 são atribuídos a esse limite.
- O limite 2 é atingido quando 60% do buffer é preenchido. Os CoS 2 e 3 são atribuídos a esse limite.
- O limite 3 é atingido quando 80% do buffer é preenchido. Os CoS 4 e 5 são atribuídos a esse limite.
- O limite 4 é atingido quando 100% do buffer é preenchido. Os CoS 6 e 7 são atribuídos a esse limite.

No diagrama na [Figura 1](#), todos os pacotes com um CoS de 0 ou 1 serão descartados se o buffer estiver 50% preenchido. Todos os pacotes com um CoS de 0, 1, 2 ou 3 serão descartados se os buffers estiverem 60% preenchidos. Pacotes com um CoS 6 ou 7 são descartados quando os buffers estiverem completamente cheios.

Figure 1



Observação: assim que o preenchimento do buffer cair abaixo de um determinado limite, os pacotes com o CoS associado não serão mais descartados.

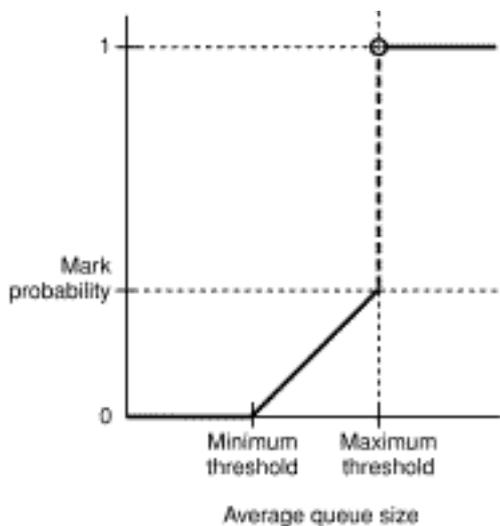
Random Early Detection e Weighted Random Early Detection

A detecção antecipada aleatória ponderada (WRED - Weighted Random Early Detection) é um mecanismo de prevenção de congestionamento que descarta aleatoriamente pacotes com uma certa precedência de IP quando os buffers atingem um limite de preenchimento definido. O WRED é uma combinação destes dois recursos:

- Queda traseira
- RED

O RED não tem conhecimento de precedência ou é compatível com CoS. O RED usa um dos limiares únicos quando o valor limite do buffer é preenchido. O RED começa a descartar aleatoriamente pacotes (mas não todos os pacotes, como na queda traseira) até que o limite máximo (máximo) seja atingido. Depois que o limite máximo é atingido, todos os pacotes são descartados. A probabilidade de um pacote ser descartado aumenta linearmente com o aumento do preenchimento do buffer acima do limite. O diagrama na [Figura 2](#) mostra a probabilidade de queda de pacote:

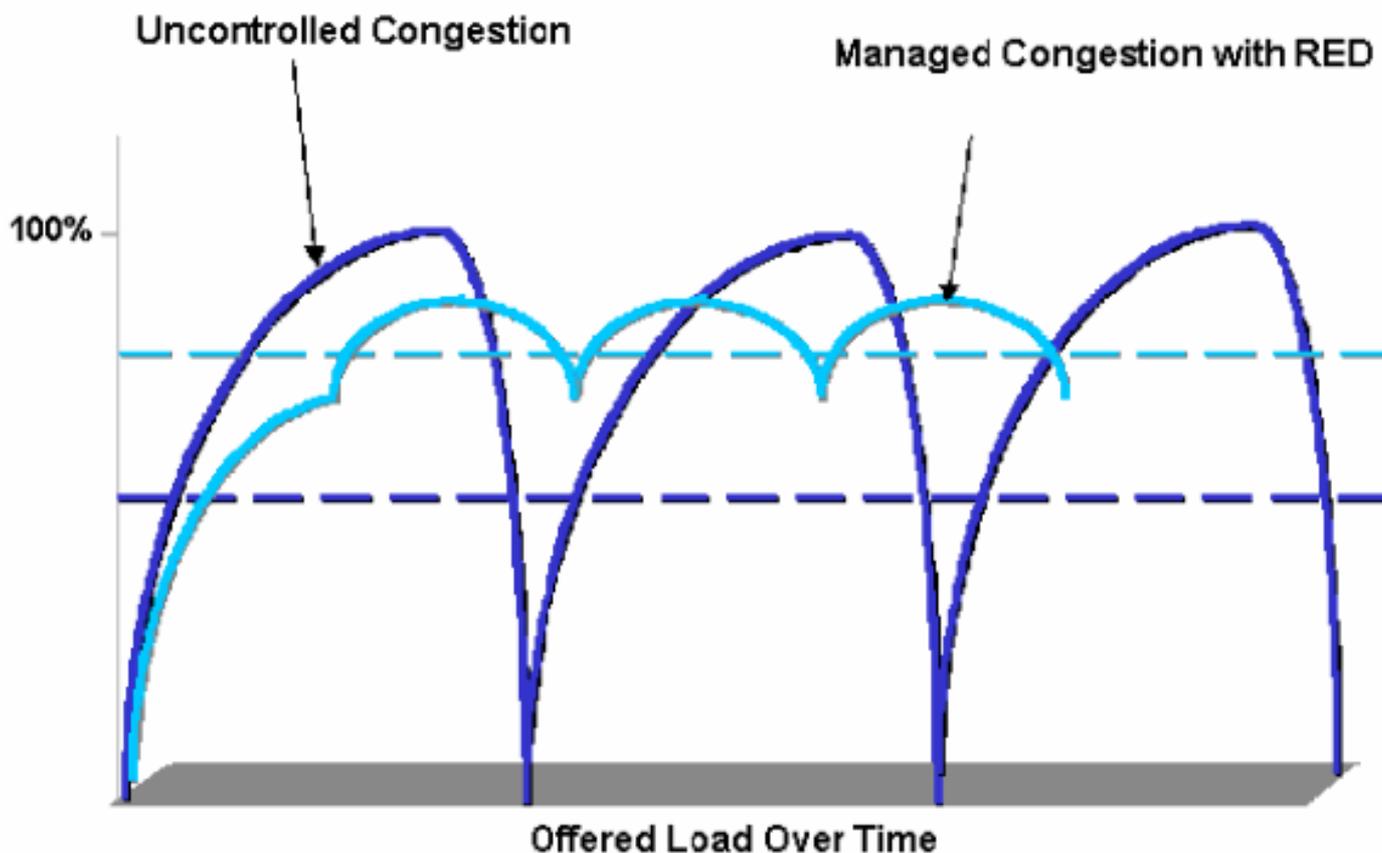
Figura 2: Probabilidade de descarte de pacote



Observação: a probabilidade da marca neste diagrama é ajustável em VERMELHO, o que significa que o declive da probabilidade de queda linear é ajustável.

RED e WRED são mecanismos muito úteis para evitar congestionamento para tráfego baseado em TCP. Para outros tipos de tráfego específicos, o RED não é muito eficiente. Isso ocorre porque o RED aproveita o mecanismo de janelamento que o TCP usa para gerenciar o congestionamento. O RED evita o congestionamento típico que ocorre em um roteador quando várias sessões TCP passam pela mesma porta do roteador. O mecanismo é chamado de sincronização de rede global. O diagrama na [Figura 3](#) mostra como o RED tem um efeito suave na carga:

Figura 3: VERMELHO para prevenção de congestionamento



Para obter mais informações sobre como o RED pode reduzir o congestionamento e suavizar o tráfego através do roteador, consulte a seção [Como o Roteador Interage com o TCP do documento Visão geral da prevenção de congestionamento](#).

O WRED é semelhante ao RED, pois ambos definem alguns limiares mínimos e, quando esses limiares mínimos são atingidos, os pacotes são descartados aleatoriamente. O WRED também define certos limites máximos e, quando esses limites máximos são atingidos, todos os pacotes são descartados. WRED também é compatível com CoS, o que significa que um ou mais valores de CoS são adicionados a cada par limiar mínimo/máximo. Quando o limiar mínimo é excedido, os pacotes são descartados aleatoriamente com o CoS atribuído. Considere este exemplo com dois limiares na fila:

- Os CoS 0 e 1 são atribuídos ao limiar mínimo 1 e ao limite máximo 1. O limiar mínimo 1 é definido como 50% do preenchimento do buffer e o limite máximo 1 é definido como 80%.
- Os CoS 2 e 3 são atribuídos ao limiar mínimo 2 e ao limite máximo 2. O limiar mínimo 2 é definido como 70% do preenchimento do buffer, e o limite máximo 2 é definido como 100%.

Assim que o buffer exceder o limite mínimo 1 (50%), os pacotes com CoS 0 e 1 começarão a ser descartados aleatoriamente. Mais pacotes são descartados à medida que a utilização do buffer cresce. Se o limiar mínimo 2 (70%) for atingido, os pacotes com CoS 2 e 3 começarão a ser descartados aleatoriamente.

Observação: neste estágio, a probabilidade de queda para pacotes com CoS 0 e 1 é muito maior que a probabilidade de queda para pacotes com CoS 2 ou CoS 3.

Sempre que o limite máximo 2 é atingido, os pacotes com CoS 0 e 1 são todos descartados, enquanto os pacotes com CoS 2 e 3 continuam a ser descartados aleatoriamente. Finalmente, quando 100% é alcançado (limite máximo 2), todos os pacotes com CoS 2 e 3 são descartados.

Os diagramas na [Figura 4](#) e na [Figura 5](#) ilustram um exemplo desses limiares:

Figura 4 - WRED com dois conjuntos de limites mínimos e limites máximos (dois serviços)

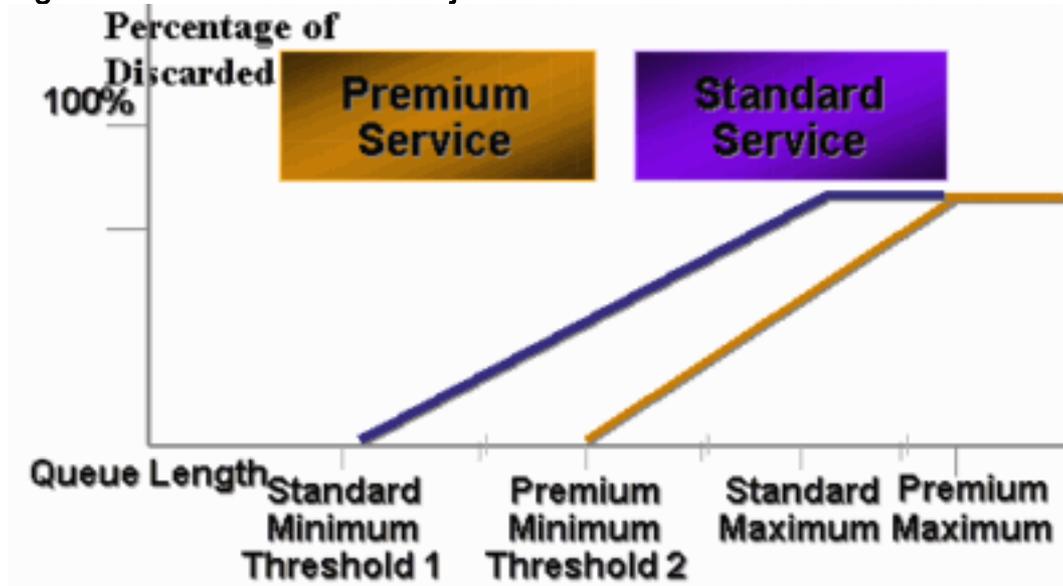
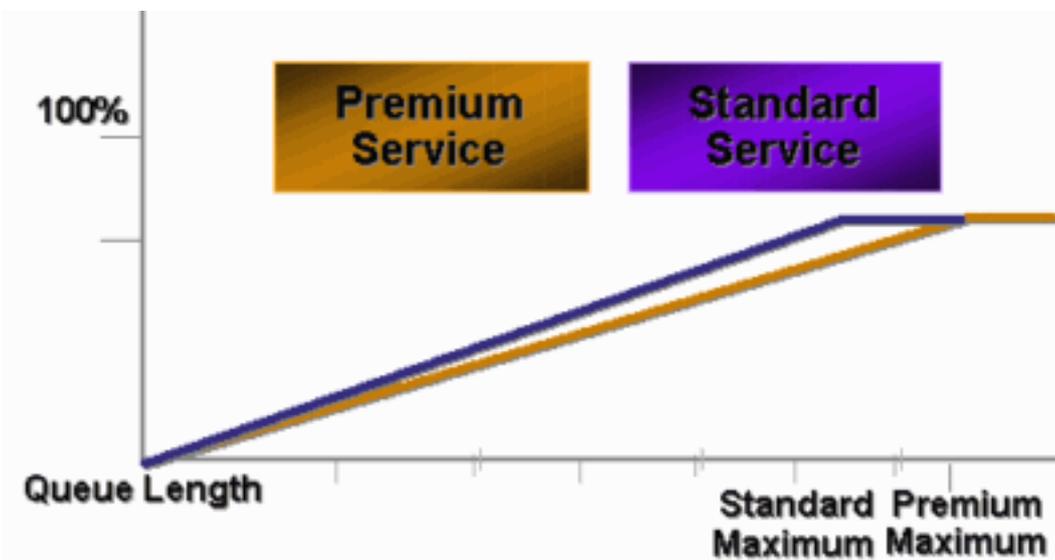


Figura 5: WRED com dois conjuntos de serviços, mas ambos os limiares mínimos são iguais a 0



A implementação antecipada do CatOS do WRED apenas definiu o limite máximo, enquanto o limite mínimo foi codificado para 0%. A parte inferior do diagrama na [Figura 5](#) destaca o comportamento resultante.

Observação: a probabilidade de descarte de um pacote é sempre não nula porque essa probabilidade está sempre acima do limite mínimo. Esse comportamento foi corrigido no software versão 6.2 e posterior.

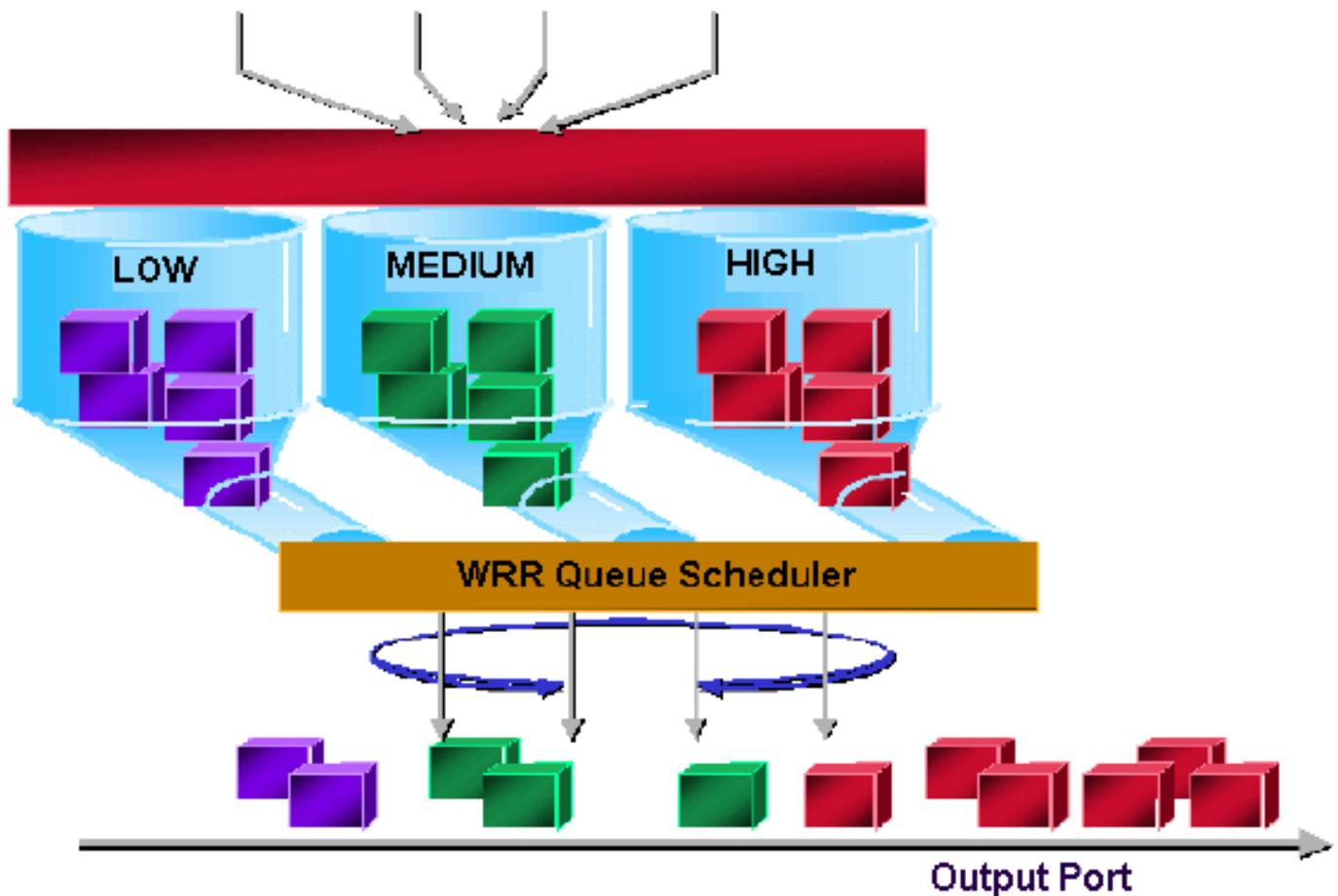
Rodízio ponderado

O rodízio ponderado (WRR) é outro mecanismo de programação de saída no Catalyst 6000. O WRR funciona entre duas ou mais filas. As filas para WRR são esvaziadas de forma round-robin e você pode configurar o peso para cada fila. Por padrão, as portas têm duas filas WRR no Catalyst 6000. O padrão é:

- Para atender à fila WRR de alta prioridade 70% do tempo
- Para atender à fila WRR de baixa prioridade 30% do tempo

O diagrama na [Figura 6](#) mostra um WRR que tem três filas servidas de forma WRR. A fila de alta prioridade (pacotes vermelhos) envia mais pacotes do que as duas outras filas:

Figura 6: Programação de saída: WRR



Observação: a maioria das 6500 placas de linha implementam WRR por largura de banda. Essa implementação de WRR por largura de banda significa que toda vez que o agendador permite que uma fila transmita pacotes, um determinado número de bytes pode ser transmitido. Esse número de bytes pode representar mais de um pacote. Por exemplo, se você enviar 5120 bytes em uma vez, poderá enviar três pacotes de 1518 bytes, para um total de 4554 bytes. Os bytes em excesso são perdidos ($5120 - 4554 = 566$ bytes). Portanto, com um peso extremo (como 1% para a fila 1 e 99% para a fila 2), o peso exato configurado pode não ser alcançado. Essa falha em atingir o peso exato é frequentemente o caso de pacotes maiores.

Algumas placas de linha de nova geração, como a 6548-RJ-45, superam essa limitação através da implementação de um rodízio ponderado pelo déficit (DWRR). O DWRR transmite das filas, mas não enfrenta a fila de baixa prioridade. O DWRR rastreia a fila de baixa prioridade que está sob transmissão e compensa na próxima rodada.

[Fila de prioridade estrita](#)

Outro tipo de fila no Catalyst 6000, uma fila de prioridade estrita, é sempre esvaziado primeiro. Assim que há um pacote na fila de prioridade estrita, o pacote é enviado.

As filas WRR ou WRED são verificadas somente depois que a fila de prioridade estrita é esvaziada. Depois que cada pacote é transmitido da fila WRR ou da fila WRED, a fila de prioridade estrita é verificada e esvaziada, se necessário.

Observação: todas as placas de linha com um tipo de enfileiramento semelhante a 1p2q1t, 1p3q8t e 1p7q8t usam DWRR. Outras placas de linha usam WRR padrão.

[Recurso de enfileiramento de saída de diferentes placas de linha](#)

no Catalyst 6000

Recursos do comando show port

Se não tiver certeza sobre o recurso de enfileiramento de uma porta, você pode emitir o comando **show port capabilities**. Esta é a saída do comando em uma placa de linha 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
```

Essa porta tem um tipo de saída de enfileiramento chamado 2q2t.

Entender o recurso de enfileiramento de uma porta

Há vários tipos de filas disponíveis nos switches Catalyst 6500/6000. As tabelas desta seção podem ficar incompletas à medida que novas placas de linha são lançadas. Novas placas de linha podem introduzir novas combinações de enfileiramento. Para obter uma descrição atual de todos os enfileiramentos disponíveis para os módulos de switch Catalyst 6500/6000, consulte a seção *Configuração de QoS* para sua versão CatOS da [Documentação de Software do Catalyst 6500 Series](#).

Observação: o Cisco Communication Media Module (CMM) não suporta todos os recursos de QoS. Verifique as notas de versão do seu software específico para determinar os recursos suportados.

Esta tabela explica a notação da arquitetura de QoS da porta:

ide Tx ¹ / Rx ²	Notação da fila	Número de filas	priorit y queue	Nº de filas WRR	Nº e tipo de limite para filas WRR
Tx	2q2t	2	—	2	2 queda traseira configurável
Tx	1p2q2t	3	1	2	2 WRED configurável
Tx	1p3q1t	4	1	3	1 WRED configurável

Tx	1p2q1t	3	1	2	1 WRED configurável
Rx	1q4t	1	—	1	4 queda traseira configurável
Rx	1p1q4t	2	1	1	4 queda traseira configurável
Rx	1p1q0t	2	1	1	Não configurável
Rx	1p1q8t	2	1	1	8 WRED configurável

¹ Tx = transmissão.

² Rx = recepção.

Esta tabela lista todos os módulos e os tipos de fila no lado Rx e Tx da interface ou porta:

Módulo	Filas Rx	Filas 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-X6224-100FX-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
------------------	------	------

[Criar QoS no Catalyst 6500/6000](#)

Três campos no Catalyst 6500/6000 são usados para fazer QoS:

- A precedência IP—Os três primeiros bits do campo Tipo de Serviço (ToS) no cabeçalho IP
- O ponto de código de serviços diferenciados (DSCP)—Os primeiros seis bits do campo ToS no cabeçalho IP
- O CoS—Os três bits usados no nível da Camada 2 (L2)Esses três bits são parte do cabeçalho Inter-Switch Link (ISL) ou estão dentro da marca IEEE 802.1Q (dot1q). Não há CoS dentro de um pacote Ethernet não marcado.

[Mecanismo de programação de saída no Catalyst 6500/6000](#)

Quando um quadro é enviado do barramento de dados a ser transmitido, o CoS do pacote é o único parâmetro considerado. O pacote passa por um agendador, que escolhe a fila na qual o pacote é colocado. Portanto, lembre-se de que o agendamento de saída e todos os mecanismos discutidos neste documento são compatíveis apenas com CoS.

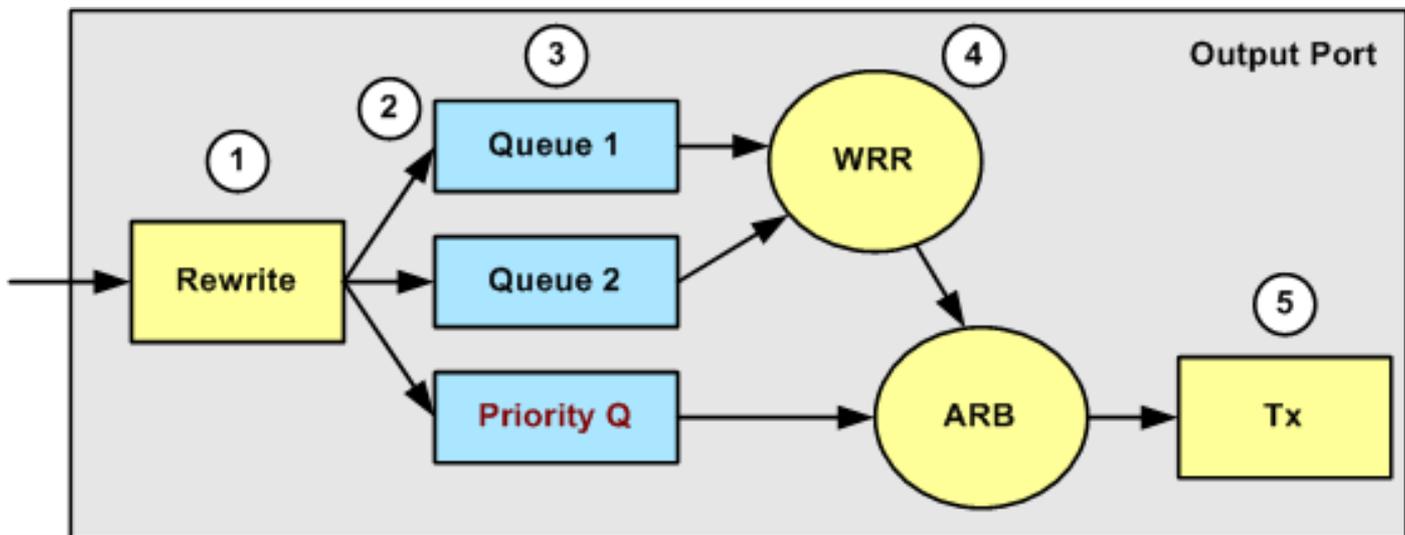
O Catalyst 6500/6000 com uma placa de recurso de switch multicamada (MSFC - Multilayer Switch Feature Card) usa um DSCP interno para classificar o pacote. O Catalyst 6500/6000 configurado com QoS habilitado atribui um valor de DSCP quando a decisão de encaminhamento é tomada no nível de PFC. Esse DSCP é atribuído a qualquer pacote, que inclui pacotes não IP, e é mapeado para o CoS para permitir o agendamento de saída. Você pode configurar o mapeamento de valores de DSCP para CoS no Catalyst 6500/6000. Se você deixar o valor padrão, poderá derivar o CoS do DSCP. A fórmula é:

DSCP_value / 8

Além disso, o valor de DSCP é mapeado no CoS do pacote de saída, se o pacote for um pacote IP marcado como ISL ou dot1q (VLAN não nativa). O valor de DSCP também é gravado no campo ToS do cabeçalho IP.

O diagrama na [Figura 7](#) mostra uma fila 1p2q2t. As filas WRR são esvaziadas com o uso do agendador WRR. Também há um árbitro que verifica entre cada pacote das filas WRR para determinar se há algo na fila de prioridade estrita.

Figura 7



1. O campo ToS é reescrito no cabeçalho IP e no campo 802.1p/ISL CoS.
2. A fila e o limite de agendamento são selecionados com base no CoS, por meio de um mapa configurável.
3. Cada fila tem limites e tamanho configuráveis, e algumas filas têm WRED.
4. A remoção da fila usa WRR entre duas filas.
5. O encapsulamento de saída pode ser dot1q, ISL ou nenhum.

[Configuração, monitoramento e programação de saída no Catalyst 6500/6000](#)

[Configuração padrão para QoS no Catalyst 6500/6000](#)

Esta seção fornece um exemplo de saída da configuração de QoS padrão em um Catalyst 6500/6000, além de informações sobre o que esses valores significam e como você pode ajustar os valores.

A QoS é desativada por padrão quando você emite este comando:

```
set qos disable
```

Os comandos nesta lista mostram a atribuição padrão para cada CoS em uma porta 2q2t. A fila 1 tem CoS 0 e 1 atribuídos ao primeiro limite e CoS 2 e 3 atribuídos ao segundo limite. A fila 2 tem CoS 4 e 5 atribuídos ao seu primeiro limiar e os CoS 6 e 7 atribuídos ao seu segundo limiar:

```
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
```

Esses comandos exibem o nível de limite por padrão em uma porta 2q2t para cada fila:

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

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

Você pode atribuir o peso padrão a cada uma das filas WRR. Emita este comando para atribuir os pesos padrão para a fila 1 e a fila 2:

Observação: a fila de baixa prioridade é atendida 5/260 por cento do tempo, e a fila de alta prioridade é atendida 255/260 por cento do tempo.

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

A disponibilidade total do buffer é dividida entre as duas filas. A fila de baixa prioridade é atribuída corretamente a 80% dos buffers disponíveis porque essa é a fila que provavelmente tem pacotes armazenados em buffer e sentados por algum tempo. Execute este comando para definir a disponibilidade:

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

Você pode visualizar configurações semelhantes para a porta 1p2q2t nesta configuração:

```
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
```

Observação: por padrão, CoS 5 (tráfego de voz) é atribuído à fila de prioridade estrita.

Configuração

A primeira etapa da configuração é ativar a QoS. Lembre-se de que a QoS está desativada, por padrão. Quando a QoS é desabilitada, o mapeamento de CoS é irrelevante. Há uma única fila que é servida como FIFO, e todos os pacotes são descartados lá.

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

```
QoS is enabled
```

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

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

O valor de CoS precisa ser atribuído à fila ou ao limite para todos os tipos de fila. O mapeamento definido para um tipo de porta 2q2t não é aplicado a nenhuma porta 1p2q2t. Além disso, o mapeamento feito para 2q2t é aplicado a todas as portas que têm um mecanismo de enfileiramento 2q2t. Emita este comando:

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

Observação: as filas são sempre numeradas para começar com a fila de prioridade mais baixa possível e terminar com a fila de prioridade estrita disponível. Aqui está um exemplo:

- A fila 1 é a fila WRR de baixa prioridade
- A fila 2 é a fila WRR de alta prioridade
- A fila 3 é a fila de prioridade estrita

Você deve repetir esta operação para todos os tipos de filas. Caso contrário, você manterá a atribuição de CoS padrão. Aqui está um exemplo para 1p2q2t:

Configuração

```
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
```

[Saída do console](#)

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

QoS tx priority queue and threshold mapped to cos successfully

Você deve configurar o peso de WRR para as duas filas de WRR. Emita este comando:

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

Weight_1 se refere à fila 1, que deve ser a fila WRR de baixa prioridade. *Weight_1* deve ser sempre menor que *weight_2*. O peso pode ter qualquer valor entre 1 e 255. Você pode atribuir a porcentagem com estas fórmulas:

- +B650Fila 1:

$$\text{weight_1} / (\text{weight_1} + \text{weight_2})$$

- Fila 2:

$$\text{weight_2} / (\text{weight_1} + \text{weight_2})$$

Você também deve definir o peso para os vários tipos de filas. O peso não precisa ser o mesmo. Por exemplo, para o 2q2t, em que a fila 1 é atendida 30 por cento do tempo e a fila 2 é atendida 70 por cento do tempo, você pode emitir este comando para definir o 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.

[Saída do console](#)

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

QoS wrr ratio is set successfully

Você também deve definir a proporção da fila de transmissão, que se refere à forma como os buffers são divididos entre as diferentes filas. Emita este comando:

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

Observação: se você tiver três filas (1p2q2t), deverá definir a fila WRR de alta prioridade e a fila de prioridade estrita no mesmo nível por razões de hardware.

[Configuração](#)

```
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.

Saída do 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

Como essa saída de console ilustra, a soma dos valores da fila deve ser 100. Deixe a maior parte dos buffers para a fila WRR de baixa prioridade porque essa fila precisa do maior buffer. As outras filas são servidas com prioridade mais elevada.

A última etapa é configurar o nível de limiar para a fila WRED ou para a fila de queda traseira. Execute estes comandos:

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

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

Configuração

```
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

Saída do 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 1p2q2t tx queue 1 50 90
```

WRED thresholds for queue 1 set to 50 and 90 on all WRED-capable 1p2q2t ports

```
tamer (enable) set qos wred 1p2q2t tx queue 2 40 80
```

WRED thresholds for queue 2 set to 40 and 80 on all WRED-capable 1p2q2t ports

O comando **set qos wred 1p2q2t tx queue 2 40 80** funciona em conjunto com o CoS para o mapeamento de limiares. Por exemplo, ao emitir os comandos na lista abaixo, você garante que—na porta 1p2q2t na direção de transmissão—os pacotes com CoS 0, 1, 2 e 3 sejam enviados na primeira fila (a fila WRR baixa). Quando os buffers nessa fila estão 50% preenchidos, o WRED começa a descartar pacotes com CoS 0 e 1. Os pacotes com CoS 2 e 3 são descartados somente quando os buffers na fila são 90% preenchidos.

```
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 wred 1p2q2t tx queue 1 50 90
```

[Monitorar o Agendamento de Saída e Verificar a Configuração](#)

Um comando simples a ser usado para verificar a configuração atual do tempo de execução para a programação de saída de uma porta é **show qos info runtime *mod/port***. O comando exibe estas informações:

- O tipo de enfileiramento na porta
- O mapeamento de CoS para as diferentes filas e limiares
- O compartilhamento do buffer
- O peso do WRR

Neste exemplo, os valores estão em 20% de WRR para a fila 1 e 80% de WRR para a fila 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 : 1p2q2t

Rx port type of port 1/1 : 1p1q4t

Interface type: port-based

ACL attached:

The qos trust type is set to untrusted

Default CoS = 0

Queue and Threshold Mapping for 1p2q2t (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)

```

No próximo exemplo, observe que os pesos WRR não são o valor padrão de 1. Os pesos foram definidos para os valores de 20 para a fila 1 e 80 para a fila 2. Este exemplo usa um gerador de tráfego para enviar 2 Gb de tráfego para um Catalyst 6000. Esses 2 Gb de tráfego devem sair pela porta 1/1. Como a porta 1/1 está com excesso de assinaturas, muitos pacotes são descartados (1 Gbps). O comando **show mac** mostra que há muita queda de saída:

```
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
```

Considere os pacotes que são descartados. É assim que o padrão de tráfego sugerido é dividido:

- 1 Gb de tráfego com precedência de IP 0
- 250 Mb de tráfego com precedência de IP 4
- 250 Mb de tráfego com precedência de IP 5
- 250 Mb de tráfego com precedência de IP 6
- 250 Mb de tráfego com precedência de IP 7

De acordo com o mapeamento de CoS, esse tráfego é enviado:

- 1 Gb de tráfego para a fila 1 limite 1
- 0 Mb de tráfego para a fila 1 limite 2
- 500 Mb de tráfego para a fila 2 limite 1
- 250 Mb de tráfego para a fila 2 limite 2
- 250 Mb de tráfego para a fila 3 (fila de prioridade estrita)

O switch deve confiar no tráfego recebido para que a precedência de IP de entrada seja preservada no switch e seja usado para mapear para o valor de CoS para programação de saída.

Observação: a precedência de IP padrão para o mapeamento de CoS é precedência de IP igual a CoS.

Emita o comando **show qos stat 1/1** para ver os pacotes que foram descartados e a porcentagem aproximada:

- Neste ponto, nenhum pacote é descartado na fila 3 (CoS 5).
- 91,85% dos pacotes descartados são pacotes CoS 0 na fila 1.
- 8% dos pacotes descartados são CoS 4 e 6 na fila 2, limiar 1.
- 0,15 por cento dos pacotes descartados são CoS 7 na fila 2, limite 2.

Esta saída ilustra o uso do comando:

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

```
Tx port type of port 1/1 : lp2q2t
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 : lplq4t
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 você alterar o peso do WRR de volta para o valor padrão depois que os contadores tiverem sido limpos, somente 1% dos pacotes descartados ocorrerão na fila 2 em vez dos 8% que apareceram anteriormente:

Observação: o valor padrão é 5 para a fila 1 e 255 para a fila 2.

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

```
TX port type of port 1/1 : lp2q2t
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 : lplq4t
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

Usar o Agendamento de Saída para Reduzir o Atraso e o Jitter

O exemplo na seção [Monitorar a Programação de Saída e Verificar a Configuração](#) demonstra o benefício da implementação de programação de saída, que evita uma queda de VoIP ou tráfego de missão crítica em caso de excesso de assinatura da porta de saída. O excesso de assinaturas ocorre com pouca frequência em uma rede normal, particularmente em um link Gigabit. Geralmente, a assinatura excessiva só ocorre durante horários de pico de tráfego ou durante surtos de tráfego em um período de tempo muito curto.

Mesmo sem excesso de assinaturas, o agendamento de saída pode ser de grande ajuda em uma rede onde a QoS é implementada de ponta a ponta. O agendamento de saída ajuda a reduzir o atraso e a instabilidade. Esta seção fornece exemplos de como o agendamento de saída pode ajudar a reduzir o atraso e a instabilidade.

Reduza o atraso

O atraso de um pacote é aumentado pelo tempo "perdido" no buffer de cada switch durante a espera pela transmissão. Por exemplo, um pequeno pacote de voz com um CoS de 5 é enviado de uma porta durante um grande backup ou transferência de arquivos. Se você não tiver nenhum QoS para a porta de saída e se assumir que o pequeno pacote de voz está enfileirado após pacotes de 10 bytes grandes, você pode facilmente calcular o tempo de velocidade Gigabit necessário para transmitir os 10 pacotes grandes:

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

Se esse pacote precisar atravessar oito ou nove switches enquanto passa pela rede, pode ocorrer um atraso de aproximadamente 1 ms. Esse valor conta apenas atrasos na fila de saída do switch que é atravessado na rede.

Observação: se você precisar colocar em fila os mesmos 10 pacotes grandes em uma interface de 10 Mbps (por exemplo, com um telefone IP e um PC conectado), o atraso que é apresentado é:

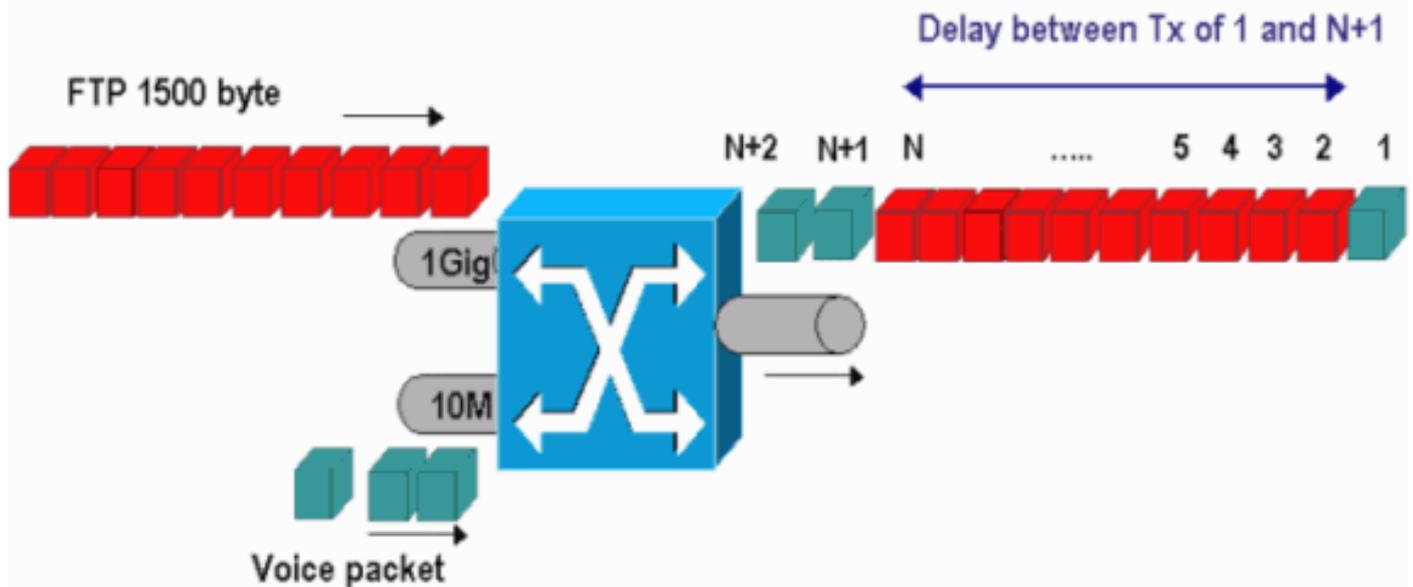
`(10 × 1500 × 8) = 120,000 bits that are transmitted in 12 ms`

A implementação de programação de saída garante que os pacotes de voz com um CoS de 5 sejam colocados na fila de prioridade estrita. Esses posicionamentos garantem que esses pacotes sejam enviados antes de qualquer pacote com um CoS inferior a 5, o que reduz os atrasos.

Reduzir o atraso de sincronismo

Outro benefício importante da implementação do agendamento de saída é que ela reduz o jitter. Jitter é a variação no atraso observada para pacotes dentro do mesmo fluxo. O diagrama na [Figura 8](#) mostra um exemplo de cenário de como o agendamento de saída pode reduzir o jitter:

Figura 8



Neste cenário, há dois fluxos que uma única porta de saída deve enviar:

- Um fluxo de voz recebido em uma porta Ethernet de 10 Mbps
- Um fluxo FTP que está sendo recebido em um uplink Ethernet de 1 Gbps

Os dois fluxos deixam o switch através da mesma porta de saída. Este exemplo mostra o que pode acontecer sem o uso do agendamento de saída. Todos os grandes pacotes de dados podem ser intercalados entre dois pacotes de voz, o que cria jitter na recepção do pacote de voz do mesmo fluxo. Há um atraso maior entre a recepção do pacote n e do pacote $n+1$ à medida que o switch transmite o pacote de dados grande. No entanto, o atraso entre $n+1$ e $n+2$ é insignificante. Isso resulta em instabilidade no fluxo de tráfego de voz. Você pode facilmente evitar esse problema com o uso de uma fila de prioridade estrita. Verifique se o valor de CoS dos pacotes de voz está mapeado para a fila de prioridade estrita.

[Informações Relacionadas](#)

- [Programação de saída de QoS em Switches Catalyst 6500/6000 Series executando o Cisco IOS System Software](#)
- [Entendendo a qualidade do serviço nos Switches da família Catalyst 6000](#)
- [Páginas de Suporte de Produtos de LAN](#)
- [Página de suporte da switching de LAN](#)
- [Suporte Técnico e Documentação - Cisco Systems](#)