

# QoS-Ausgabeplanung für Catalyst Switches der Serien 6500 und 6000 mit CatOS-Systemsoftware

## Inhalt

[Einführung](#)

[Voraussetzungen](#)

[Anforderungen](#)

[Verwendete Komponenten](#)

[Konventionen](#)

[Hintergrundinformationen](#)

[Ausgabewarteschlangen fällt aus](#)

[Warteschlangentypen, die an der Ausgabeplanung beim Catalyst 6500/6000 beteiligt sind](#)

[Tail Drop](#)

[Frühzeitige Erkennung nach Zufallsprinzip und weighted Random Early Detection](#)

[Weighted Round Robin](#)

[Warteschlange mit strikter Priorität](#)

[Ausgangswarteschlangen-Funktion verschiedener Line Cards auf dem Catalyst 6000](#)

[show port-Befehlsfunktionen](#)

[Analyse der Warteschlangenkapazität eines Ports](#)

[Erstellen von QoS auf dem Catalyst 6500/6000](#)

[Ausgabeplanmechanismus für den Catalyst 6500/6000](#)

[Konfiguration, Überwachung und Ausgabeplanung auf dem Catalyst 6500/6000](#)

[Standardkonfiguration für QoS auf dem Catalyst 6500/6000](#)

[Konfiguration](#)

[Überwachen der Ausgabeplanung und Überprüfen der Konfiguration](#)

[Verwenden der Ausgabeplanung zur Verringerung von Verzögerungen und Jitter](#)

[Verzögerung reduzieren](#)

[Jitter reduzieren](#)

[Zugehörige Informationen](#)

## **Einführung**

Die Ausgabeplanung stellt sicher, dass bei starker Überbelegung kein wichtiger Datenverkehr verworfen wird. In diesem Dokument werden alle Techniken und Algorithmen erläutert, die bei der Ausgabenplanung für Cisco Catalyst Switches der Serien 6500 und 6000 mit Catalyst OS (CatOS)-Systemsoftware zum Einsatz kommen. Dieses Dokument bietet außerdem einen kurzen Überblick über die Warteschlangenfunktionen von Catalyst 6500/6000-Switches und die Konfiguration der verschiedenen Parameter für die Ausgabeplanung.

**Hinweis:** Wenn Sie Cisco IOS® Software auf Ihrem Catalyst 6500/6000 ausführen, finden Sie weitere Informationen unter [QoS-Ausgabeplanung für Catalyst Switches der Serien 6500/6000 mit Cisco IOS-Systemsoftware](#).

## Voraussetzungen

### Anforderungen

Für dieses Dokument bestehen keine speziellen Anforderungen.

### Verwendete Komponenten

Die Beispiele in diesem Dokument wurden aus einem Catalyst 6000 mit Supervisor Engine 1A und Policy Feature Card (PFC) erstellt. Die Beispiele gelten jedoch auch für eine Supervisor Engine 2 mit PFC2 oder eine Supervisor Engine 720 mit PFC3.

Die Informationen in diesem Dokument wurden von den Geräten in einer bestimmten Laborumgebung erstellt. Alle in diesem Dokument verwendeten Geräte haben mit einer leeren (Standard-)Konfiguration begonnen. Wenn Ihr Netzwerk in Betrieb ist, stellen Sie sicher, dass Sie die potenziellen Auswirkungen eines Befehls verstehen.

### Konventionen

Weitere Informationen zu Dokumentkonventionen finden Sie unter [Cisco Technical Tips Conventions](#) (Technische Tipps zu Konventionen von Cisco).

## Hintergrundinformationen

### Ausgabewarteschlangen fällt aus

Ausgabeverwerfungen werden durch eine überlastete Schnittstelle verursacht. Eine häufige Ursache hierfür ist der Datenverkehr von einer Verbindung mit hoher Bandbreite, der zu einer Verbindung mit niedriger Bandbreite umgeschaltet wird, oder der Datenverkehr von mehreren eingehenden Verbindungen, die auf eine einzige ausgehende Verbindung umgeschaltet werden.

Wenn z. B. ein großer Teil des Datenverkehrs über eine Gigabit-Schnittstelle eingeht und auf eine 100-Mbit/s-Schnittstelle ausgeschaltet wird, kann dies dazu führen, dass die Ausgangsleistung auf der 100-Mbit/s-Schnittstelle sinken. Der Grund hierfür ist, dass die Ausgabewarteschlange an dieser Schnittstelle durch den übermäßigen Datenverkehr aufgrund der Geschwindigkeitsungleichheit zwischen der eingehenden und der ausgehenden Bandbreite überlastet wird. Die Datenverkehrsrate auf der Ausgangsschnittstelle kann nicht alle Pakete akzeptieren, die gesendet werden sollen.

Die ultimative Lösung zur Lösung des Problems ist die Erhöhung der Leitungsgeschwindigkeit. Es gibt jedoch Möglichkeiten, Ausgabeverfahren zu verhindern, zu verringern oder zu steuern, wenn Sie die Leitungsgeschwindigkeit nicht erhöhen möchten. Ausgabeverwerfungen können nur dann verhindert werden, wenn Ausgabeverwerfen auf kurze Datenspitzen zurückzuführen sind. Wenn Ausgabetrophen durch einen konstanten Hochdatenfluss verursacht werden, können Sie diese Verwerfungen nicht verhindern. Sie können diese jedoch steuern.

# Warteschlangentypen, die an der Ausgabeplanung beim Catalyst 6500/600 beteiligt sind

## Tail Drop

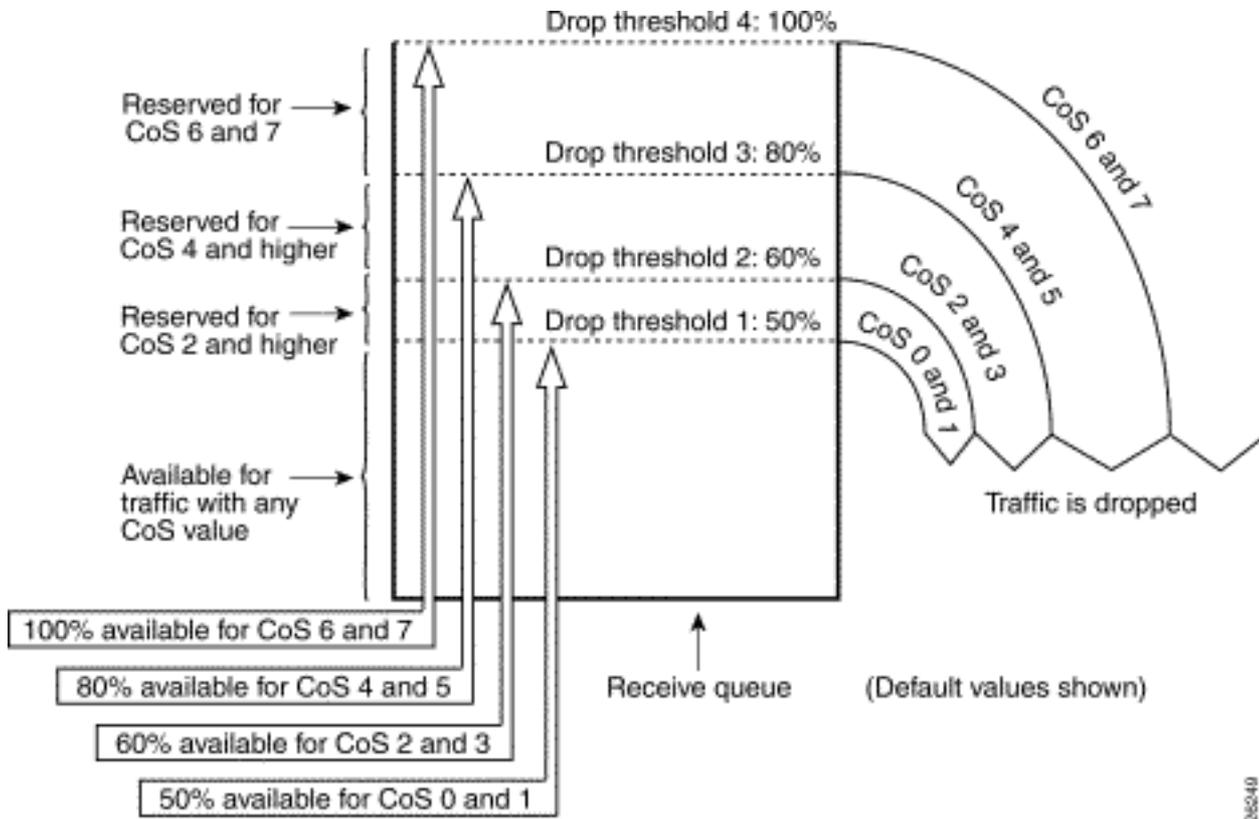
Tail Drop ist ein grundlegender Überlastungsvermeidungsmechanismus. Tail Drop behandelt den gesamten Datenverkehr gleich und unterscheidet nicht zwischen CoS (Classes of Service), wenn Warteschlangen in Zeiten von Überlastung zu füllen beginnen. Wenn die Ausgabewarteschlange voll ist und der Tail-Drop aktiv ist, werden Pakete verworfen, bis die Überlastung beseitigt ist und die Warteschlange nicht mehr voll ist. Das Tail Drop ist der einfachste Typ der Überlastungsvermeidung und berücksichtigt keine QoS-Parameter.

Der Catalyst 6000 hat eine erweiterte Version zur Vermeidung von Engpässen beim Tail-Drop implementiert, bei der alle Pakete mit einem bestimmten CoS verworfen werden, wenn ein bestimmter Prozentsatz der Pufferüberlastung erreicht wird. Bei einem gewichteten Tail-Drop können Sie eine Reihe von Schwellenwerten definieren und jedem Schwellenwert eine CoS zuordnen. Im Beispiel in diesem Abschnitt gibt es vier mögliche Schwellenwerte. Jeder Schwellenwert wird wie folgt definiert:

- Der Schwellenwert 1 wird erreicht, wenn 50 % des Puffers gefüllt sind. Diesem Schwellenwert sind CoS 0 und 1 zugewiesen.
- Der Schwellenwert 2 wird erreicht, wenn 60 % des Puffers gefüllt sind. CoS 2 und 3 werden diesem Schwellenwert zugewiesen.
- Der Schwellenwert 3 wird erreicht, wenn 80 Prozent des Puffers gefüllt sind. CoS 4 und 5 werden diesem Schwellenwert zugewiesen.
- Der Schwellenwert 4 wird erreicht, wenn 100 Prozent des Puffers gefüllt sind. CoS 6 und 7 werden diesem Schwellenwert zugewiesen.

Im Diagramm in [Abbildung 1](#) werden alle Pakete mit einer CoS von 0 oder 1 verworfen, wenn der Puffer zu 50 % gefüllt ist. Alle Pakete mit einer CoS von 0, 1, 2 oder 3 werden verworfen, wenn die Puffer zu 60 % gefüllt sind. Pakete mit einer CoS von 6 oder 7 werden verworfen, wenn die Puffer vollständig gefüllt sind.

### **Abbildung 1**



**Hinweis:** Sobald die Puffernfüllung unter einen bestimmten Schwellenwert fällt, werden Pakete mit dem zugehörigen CoS nicht mehr verworfen.

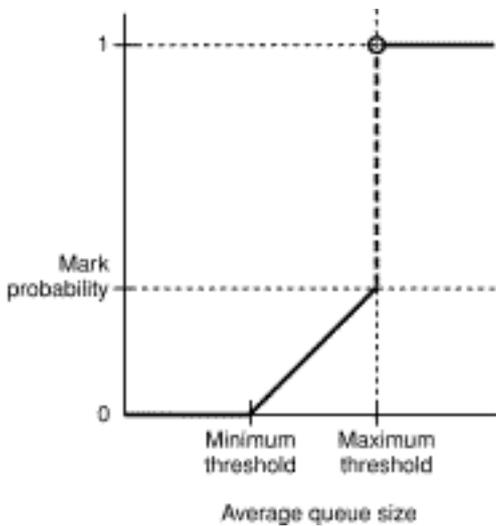
## Frühzeitige Erkennung nach Zufallsprinzip und weighted Random Early Detection

Weighted Random Early Detection (WRED) ist ein Überlastungsvermeidungsmechanismus, der Pakete mit einer bestimmten IP-Priorität willkürlich verwirft, wenn die Puffer einen festgelegten Füllungsschwellenwert erreichen. WRED ist eine Kombination dieser beiden Funktionen:

- Tail Drop
- Random Early Detection (RED)

ROT ist nicht prioritätsorientiert oder CoS-basiert. ROT verwendet einen der einzelnen Schwellenwerte, wenn der Schwellenwert für den Puffer voll ist. ROT beginnt, Pakete zufällig (aber nicht alle Pakete wie im Tail Drop) zu verworfen, bis der maximale (max.) Grenzwert erreicht ist. Wenn der maximale Grenzwert erreicht ist, werden alle Pakete verworfen. Die Wahrscheinlichkeit, dass ein Paket verworfen wird, steigt linear, wenn die Anzahl der Pufferfüllungen über den Schwellenwert steigt. Das Diagramm in [Abbildung 2](#) zeigt die Wahrscheinlichkeit von Paketverlusten:

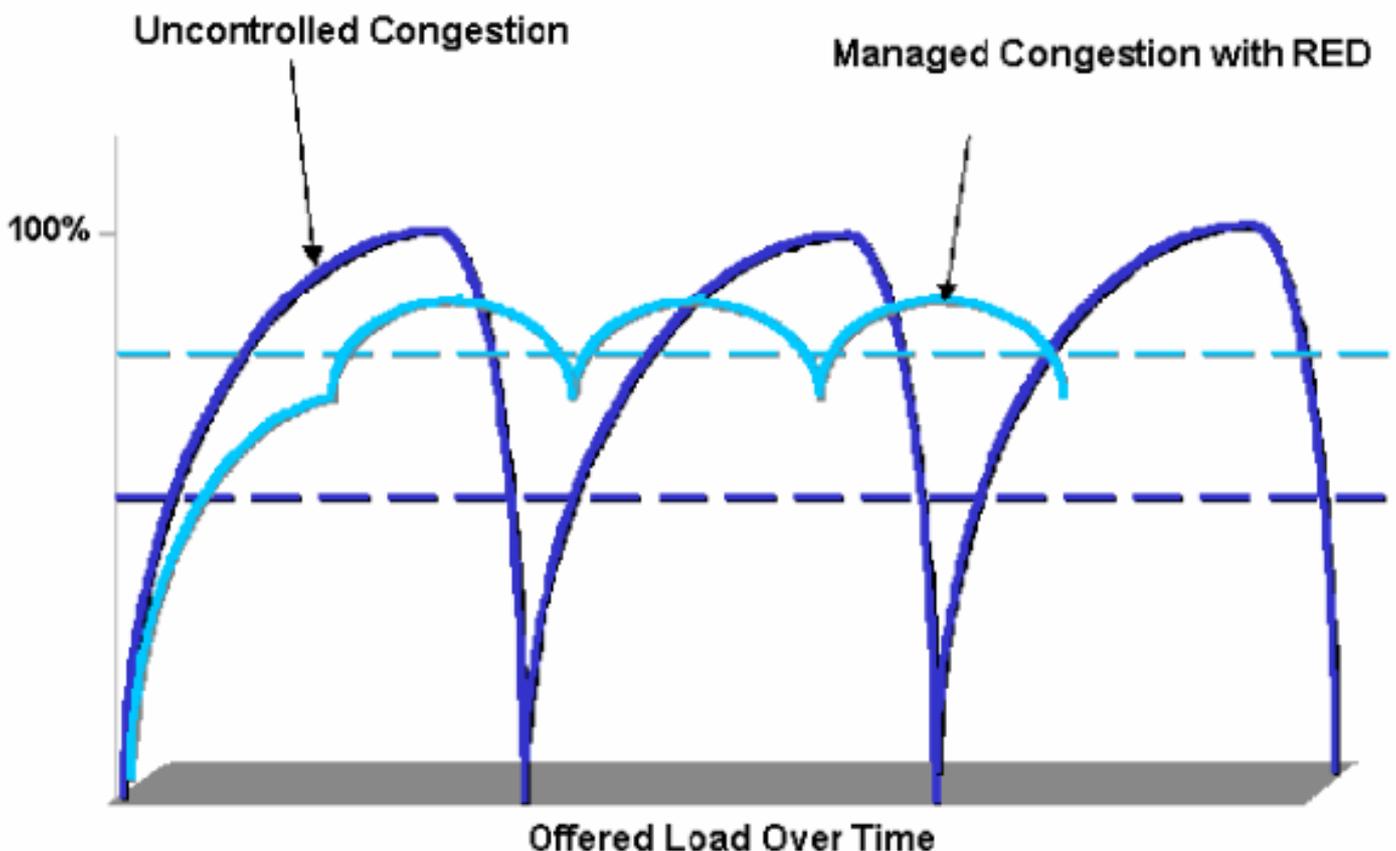
**Abbildung 2: Wahrscheinlichkeit der Paketverwerfen**



**Hinweis:** Die Markierungswahrscheinlichkeit in diesem Diagramm ist ROT einstellbar, d. h. die Neigung der linearen Fallwahrscheinlichkeit ist einstellbar.

ROT und WRED sind sehr nützliche Überlastungsvermeidungsmechanismen für TCP-basierten Datenverkehr. Bei anderen Verkehrsarten ist ROT nicht sehr effizient. ROTE nutzt den Fenstermechanismus, den TCP zur Überlastungsverwaltung verwendet. ROT vermeidet die typische Überlastung eines Routers, wenn mehrere TCP-Sitzungen über denselben Router-Port laufen. Dieser Mechanismus wird als globale Netzwerksynchronisierung bezeichnet. Das Diagramm in [Abbildung 3](#) zeigt, wie ROT eine Glättung der Last bewirkt:

Abbildung 3: ROT zur Vermeidung von Überlastungen



Weitere Informationen darüber, wie ROT Engpässe reduzieren und den Datenverkehr durch den Router glätten kann, finden Sie im [Abschnitt \*Wie der Router mit TCP interagiert im Dokument Übersicht zur Überlastungsvermeidung\*](#).

WRED ähnelt ROT insofern, als sowohl Mindestschwellenwerte (min.) als auch Pakete willkürlich verworfen werden, wenn diese Mindestschwellenwerte erreicht werden. WRED definiert außerdem bestimmte maximale Schwellenwerte und, wenn diese maximalen Schwellenwerte erreicht werden, werden alle Pakete verworfen. WRED ist auch CoS-basiert, d. h., dass jedem Min.-Schwellenwert/Max-Schwellenwertpaar mindestens ein CoS-Wert hinzugefügt wird. Wenn der Mindestschwellenwert überschritten wird, werden Pakete willkürlich verworfen, wobei der CoS zugewiesen wird. Betrachten Sie dieses Beispiel mit zwei Schwellenwerten in der Warteschlange:

- CoS 0 und 1 werden dem Min.-Schwellenwert 1 und dem Höchstwert 1 zugewiesen. Der Mindestwert 1 ist auf 50 % der Pufferbelegung und der Höchstwert 1 auf 80 % festgelegt.
- CoS 2 und 3 werden dem Min.-Schwellenwert 2 und dem Mindest-Schwellenwert 2 zugewiesen. Der Mindestwert 2 ist auf 70 % der Pufferbelegung und der Höchstwert 2 auf 100 % festgelegt.

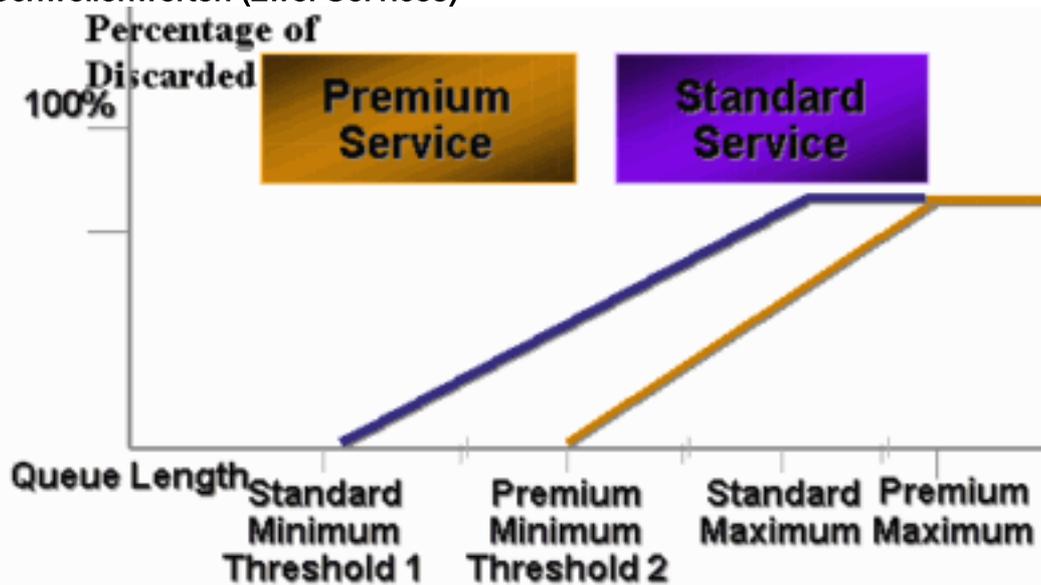
Sobald der Puffer den Min.-Schwellenwert 1 (50 %) überschreitet, werden Pakete mit CoS 0 und 1 willkürlich verworfen. Bei wachsender Puffer-Auslastung werden mehr Pakete verworfen. Wenn der Mindestwert 2 (70 %) erreicht ist, werden Pakete mit CoS 2 und 3 nach dem Zufallsprinzip verworfen.

**Hinweis:** In dieser Phase ist die Verlustrate bei Paketen mit CoS 0 und 1 bei Paketen mit CoS 2 oder CoS 3 deutlich höher als die Verlustrate.

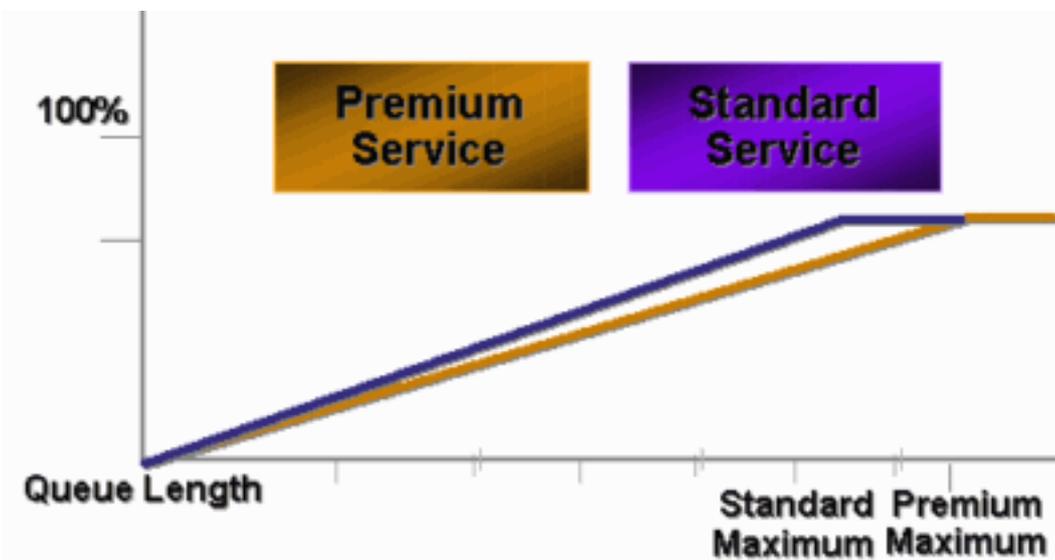
Wenn der maximale Grenzwert 2 erreicht wird, werden alle Pakete mit CoS 0 und 1 verworfen, während Pakete mit CoS 2 und 3 weiterhin zufällig verworfen werden. Wenn schließlich 100 Prozent erreicht sind (max. Grenzwert 2), werden alle Pakete mit CoS 2 und 3 verworfen.

Die Diagramme in [Abbildung 4](#) und [Abbildung 5](#) veranschaulichen ein Beispiel für folgende Schwellenwerte:

**Abbildung 4: WRED mit zwei Gruppen von Mindestschwellenwerten und maximalen Schwellenwerten (zwei Services)**



**Abbildung 5: WRED mit zwei Servicesätzen, aber beide Min.-Schwellenwerte sind 0**



Die frühe CatOS-Implementierung von WRED setzte nur den maximalen Schwellenwert fest, während der Mindestwert fest auf 0 % festgelegt war. Im unteren Teil des Diagramms in [Abbildung 5](#) wird das resultierende Verhalten hervorgehoben.

**Hinweis:** Die Drop-Wahrscheinlichkeit für ein Paket ist immer nicht NULL, da diese Wahrscheinlichkeit immer über dem Mindestwert liegt. Dieses Verhalten wurde in der Softwareversion 6.2 und höher korrigiert.

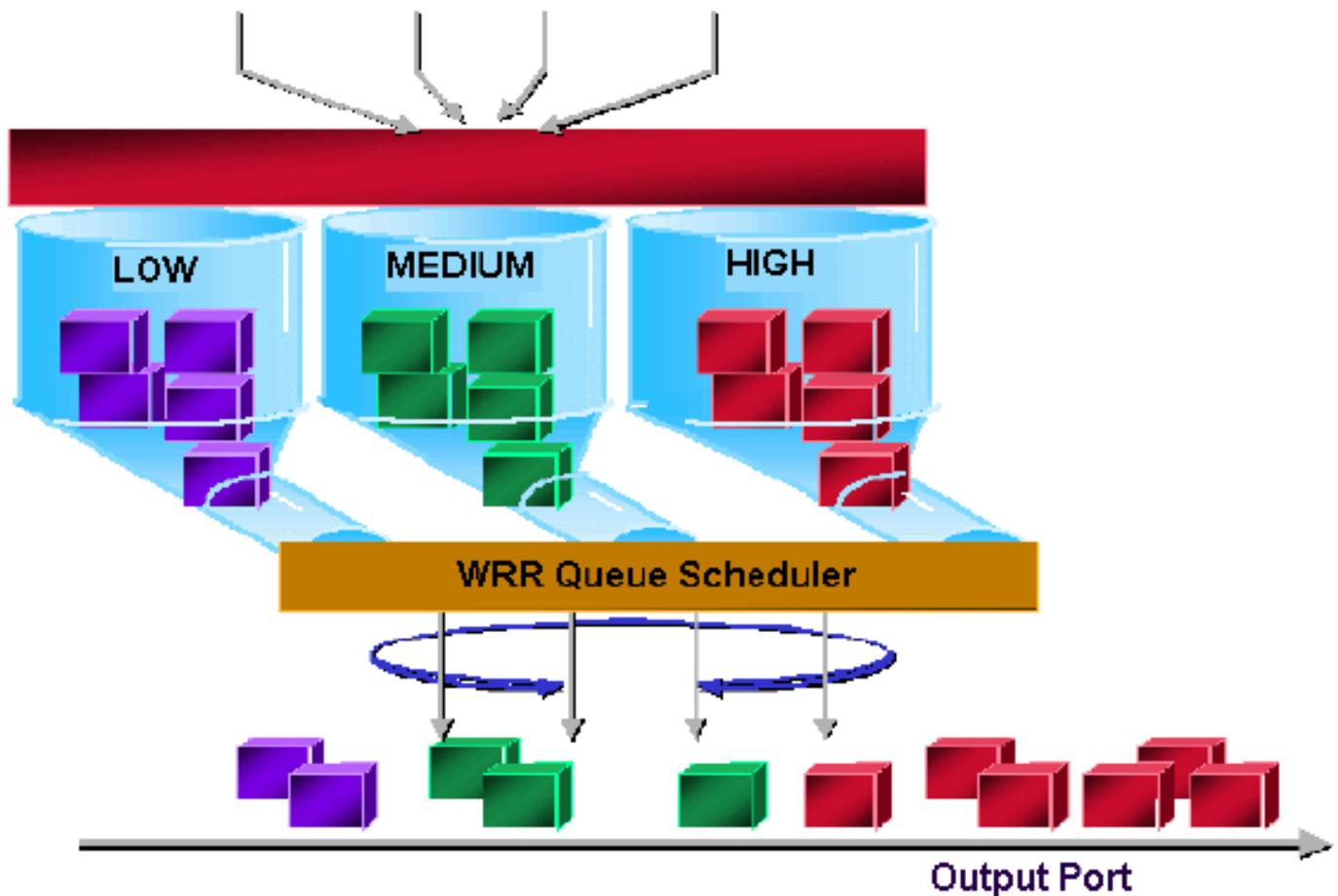
## [Weighted Round Robin](#)

Weighted Round Robin (WRR) ist ein weiterer Mechanismus für die Ausgabeplanung auf dem Catalyst 6000. WRR funktioniert zwischen zwei oder mehr Warteschlangen. Die Warteschlangen für den WRR werden rundherum geleert, und Sie können das Gewicht für jede Warteschlange konfigurieren. Standardmäßig verfügen Ports auf dem Catalyst 6000 über zwei WRR-Warteschlangen. Der Standardwert ist:

- Bereitstellung der WRR-Warteschlange mit hoher Priorität zu 70 Prozent der Zeit
- Bereitstellung der WRR-Warteschlange mit niedriger Priorität (30 Prozent der Zeit)

Das Diagramm in [Abbildung 6](#) zeigt einen WRR mit drei Warteschlangen, die auf WRR-Weise bedient werden. Die Warteschlange mit hoher Priorität (rote Pakete) sendet mehr Pakete als die beiden anderen Warteschlangen:

**Abbildung 6: Ausgabeplanung: WRR**



**Hinweis:** Die meisten Linecards der Serie 6500 implementieren WRR pro Bandbreite. Diese Implementierung von WRR pro Bandbreite bedeutet, dass jedes Mal, wenn der Scheduler eine Warteschlange zum Übertragen von Paketen zulässt, eine bestimmte Anzahl von Byte übertragen werden darf. Diese Byteanzahl kann mehr als ein Paket darstellen. Wenn Sie z. B. 5120 Byte auf einmal senden, können Sie drei 1518-Byte-Pakete senden, d. h. insgesamt 4554 Byte. Die überzähligen Bytes gehen verloren ( $5120 - 4554 = 566$  Byte). Daher wird bei einem extremen Gewicht (z. B. 1 Prozent für die Warteschlange 1 und 99 Prozent für die Warteschlange 2) das exakte konfigurierte Gewicht möglicherweise nicht erreicht. Bei größeren Paketen ist dies häufig der Fall, wenn das genaue Gewicht nicht erreicht wird.

Einige Line Cards der neuen Generation, wie der 6548-RJ-45, überwinden diese Einschränkung durch die Implementierung von Deficit Weighted Round Robin (DWRR). Der DWRR wird aus den Warteschlangen übertragen, jedoch nicht aus der Warteschlange mit niedriger Priorität. DWRR verfolgt die Warteschlange mit niedriger Priorität, die derzeit übertragen wird, und kompensiert diese in der nächsten Runde.

### Warteschlange mit strikter Priorität

Ein anderer Warteschlangentyp im Catalyst 6000, eine Warteschlange mit strikter Priorität, wird immer zuerst geleert. Sobald sich ein Paket in der Warteschlange mit strikter Priorität befindet, wird das Paket gesendet.

Die WRR- oder WRED-Warteschlangen werden erst nach Leerung der Warteschlange mit strikter Priorität überprüft. Nachdem jedes Paket entweder von der WRR-Warteschlange oder der WRED-Warteschlange übertragen wurde, wird die Warteschlange mit der strikten Priorität überprüft und gegebenenfalls geleert.

**Hinweis:** Alle Linecards mit einem Warteschlangentyp, der 1p2q1t, 1p3q8t und 1p7q8t ähnelt, verwenden DWRR. Andere Linecards verwenden Standard-WRR.

## Ausgangswarteschlangen-Funktion verschiedener Line Cards auf dem Catalyst 6000

### show port-Befehlsfunktionen

Wenn Sie sich nicht sicher sind, ob die Warteschlangenfunktion eines Ports aktiviert ist, können Sie den Befehl **show port functions (Portfunktionen anzeigen)** ausführen. Dies ist die Ausgabe des Befehls auf einer WS-X6408-GBIC-Linecard:

```

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
  
```

Dieser Port verfügt über eine Warteschlangenausgabe mit der Bezeichnung 2q2t.

### Analyse der Warteschlangenkapazität eines Ports

Für Catalyst 6500/6000-Switches sind mehrere Warteschlangentypen verfügbar. Die Tabellen in diesem Abschnitt sind möglicherweise unvollständig, wenn neue Linecards veröffentlicht werden. Neue Linecards können neue Warteschlangenkombinationen einführen. Eine aktuelle Beschreibung aller Warteschlangen, die für Catalyst 6500/6000-Switch-Module verfügbar sind, finden Sie im Abschnitt *QoS-Konfiguration* für Ihre CatOS-Version der [Catalyst 6500-Softwaredokumentation](#).

**Hinweis:** Das Cisco Communication Media Module (CMM) unterstützt nicht alle QoS-Funktionen. In den Versionshinweisen für Ihre spezifische Softwareversion finden Sie die unterstützten Funktionen.

In dieser Tabelle wird die Notation der Port-QoS-Architektur erläutert:

T x 1/ R	Warteschla ngenerkenn ung	Anzahl Warte schlan gen	Prioritäts warteschl ange	Anzahl der WRR- Wartesch	Anzahl und Schwelle nwert für
-------------------	---------------------------------	----------------------------------	---------------------------------	-----------------------------------	----------------------------------------

X <sup>2</sup> -Seite				langen	WRR-Warteschlangen
T <sub>x</sub>	2q2t	2	—	2	2 konfigurierbare Tail Drop
T <sub>x</sub>	1p2q2t	1	1	2	2 konfigurierbares WRED
T <sub>x</sub>	1p3q1t	4	1	1	1 konfigurierbares WRED
T <sub>x</sub>	1p2q1t	1	1	2	1 konfigurierbares WRED
R <sub>x</sub>	1q4t	1	—	1	4 konfigurierbare Tail Drop
R <sub>x</sub>	1p1q4t	2	1	1	4 konfigurierbare Tail Drop
R <sub>x</sub>	1p1q0t	2	1	1	Nicht konfigurierbar
R <sub>x</sub>	1p1q8t	2	1	1	8 konfigurierbare WRED

<sup>1</sup> Tx = übertragen.

<sup>2</sup> Rx = Empfangen.

In dieser Tabelle sind alle Module und Warteschlangentypen auf der Rx- und Tx-Seite der Schnittstelle oder des Ports aufgeführt:

Modul	Rx-Warteschlangen	Tx-Warteschlangen
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

## [Erstellen von QoS auf dem Catalyst 6500/6000](#)

Für die Erstellung von QoS werden drei Felder auf dem Catalyst 6500/6000 verwendet:

- Die IP-Rangfolge - Die ersten drei Bit des Felds Type of Service (ToS) im IP-Header
- Differentiated Services Code Point (DSCP) - Die ersten sechs Bit des ToS-Felds im IP-Header
- CoS - Die drei Bits, die auf Layer 2 (L2) verwendet werden. Diese drei Bits sind entweder Teil des ISL-Headers (Inter-Switch Link) oder innerhalb des IEEE 802.1Q (dot1q)-Tags. In einem nicht gekennzeichneten Ethernet-Paket befindet sich kein CoS.

## [Ausgabeplanmechanismus für den Catalyst 6500/6000](#)

Wenn ein Frame vom zu übertragenden Datenbus gesendet wird, ist die CoS des Pakets der einzige Parameter, der berücksichtigt wird. Das Paket durchläuft dann einen Scheduler, der die Warteschlange auswählt, in die das Paket eingefügt wird. Beachten Sie daher, dass die Ausgabeplanung und alle in diesem Dokument behandelten Mechanismen nur CoS-basiert sind.

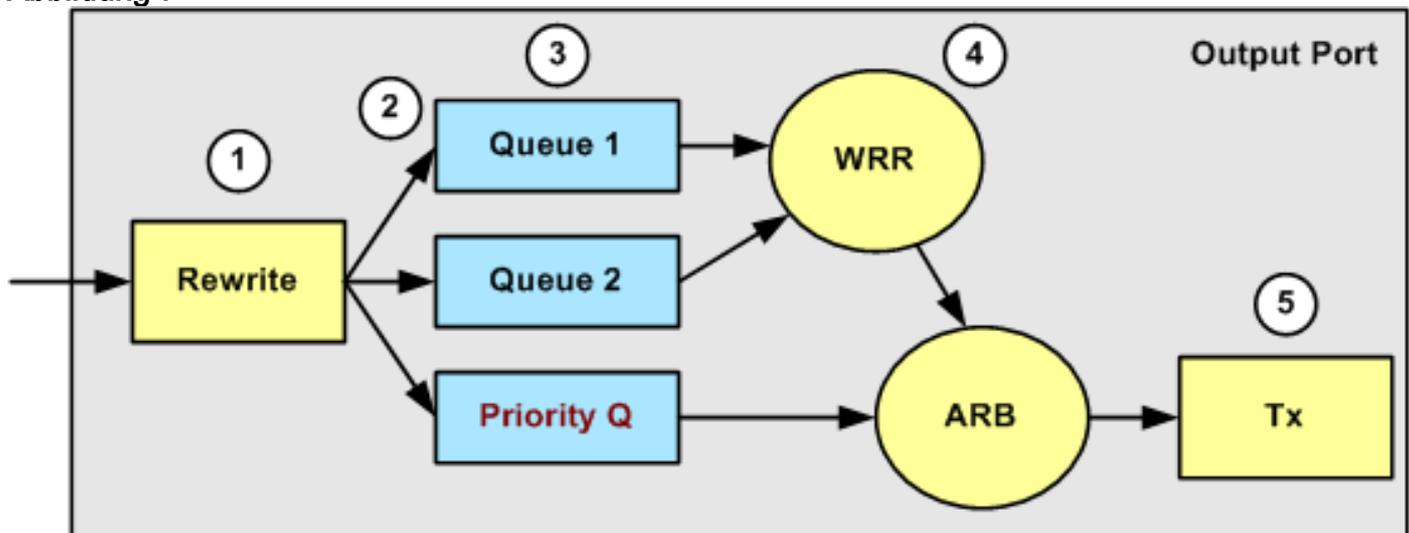
Der Catalyst 6500/6000 mit einer Multilayer Switch Feature Card (MSFC) verwendet zur Klassifizierung des Pakets ein internes DSCP. Der Catalyst 6500/6000, der mit aktivierter QoS konfiguriert wird, weist einen DSCP-Wert zu, wenn die Weiterleitungsentscheidung auf PFC-Ebene getroffen wird. Dieses DSCP wird jedem Paket zugewiesen, das auch Nicht-IP-Pakete enthält, und dem CoS zugeordnet, um die Ausgabeplanung zu ermöglichen. Sie können die Zuordnung von DSCP zu CoS-Werten auf dem Catalyst 6500/6000 konfigurieren. Wenn Sie den Standardwert beibehalten, können Sie die CoS vom DSCP ableiten. Die Formel lautet:

$DSCP\_value / 8$

Darüber hinaus wird der DSCP-Wert der CoS des ausgehenden Pakets zugeordnet, wenn es sich bei dem Paket um ein IP-Paket handelt, das mit ISL oder dot1q (nicht natives VLAN) gekennzeichnet ist. Der DSCP-Wert wird auch im ToS-Feld des IP-Headers geschrieben.

Das Diagramm in [Abbildung 7](#) zeigt eine 1p2q2t-Warteschlange. Die WRR-Warteschlangen werden mithilfe des WRR Scheduler geleert. Es gibt auch einen Schiedsrichter, der zwischen jedem Paket aus den WRR-Warteschlangen überprüft, um festzustellen, ob sich etwas in der Warteschlange mit strikter Priorität befindet.

Abbildung 7



1. Das ToS-Feld wird im IP-Header und im 802.1p/ISL CoS-Feld neu geschrieben.
2. Die Planungswarteschlange und der Schwellenwert werden anhand der CoS über eine konfigurierbare Karte ausgewählt.
3. Jede Warteschlange verfügt über konfigurierbare Größen und Schwellenwerte, und einige Warteschlangen haben WRED.
4. Beim Dequeueing wird WRR zwischen zwei Warteschlangen verwendet.
5. Die ausgehende Kapselung kann dot1q, ISL oder none sein.

## [Konfiguration, Überwachung und Ausgabeplanung auf dem Catalyst 6500/6000](#)

## Standardkonfiguration für QoS auf dem Catalyst 6500/6000

Dieser Abschnitt enthält Beispielausgaben aus der standardmäßigen QoS-Konfiguration eines Catalyst 6500/6000 sowie Informationen darüber, was diese Werte bedeuten und wie Sie die Werte einstellen können.

QoS ist standardmäßig deaktiviert, wenn Sie diesen Befehl ausführen:

```
set qos disable
```

Die Befehle in dieser Liste zeigen die Standardzuweisung für jede CoS in einem 2q2t-Port an. Für Warteschlange 1 ist CoS 0 und 1 dem ersten Schwellenwert zugewiesen, und CoS 2 und 3 ist dem zweiten Schwellenwert zugewiesen. Für Warteschlange 2 ist CoS 4 und 5 dem ersten Schwellenwert zugewiesen, und CoS 6 und 7 ist dem zweiten Schwellenwert zugewiesen:

```
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
```

Diese Befehle zeigen standardmäßig den Schwellenwert für jeden 2q2t-Port einer Warteschlange an:

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

Sie können jeder WRR-Warteschlange das Standardgewicht zuweisen. Geben Sie diesen Befehl ein, um die Standardgewichte für die Warteschlange 1 und 2 zuzuweisen:

**Hinweis:** Die Warteschlange mit niedriger Priorität wird 5/260 % der Zeit und die Warteschlange mit hoher Priorität 255/260 % der Zeit serviert.

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

Die Gesamtverfügbarkeit des Puffers wird auf die beiden Warteschlangen aufgeteilt. Die Warteschlange mit niedriger Priorität wird korrekt 80 Prozent der verfügbaren Puffer zugewiesen,

da dies die Warteschlange ist, in der Pakete höchstwahrscheinlich für einige Zeit gepuffert und gespeichert werden. Geben Sie diesen Befehl ein, um die Verfügbarkeit zu definieren:

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

Sie können ähnliche Einstellungen für den 1p2q2t-Port in dieser Konfiguration anzeigen:

```
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
```

**Hinweis:** CoS 5 (Sprachverkehr) wird standardmäßig der Warteschlange mit strikter Priorität zugewiesen.

## [Konfiguration](#)

Der erste Konfigurationsschritt ist die Aktivierung von QoS. Beachten Sie, dass QoS standardmäßig deaktiviert ist. Wenn QoS deaktiviert ist, ist die CoS-Zuordnung irrelevant. Es gibt eine einzige Warteschlange, die als FIFO bereitgestellt wird, und alle Pakete werden dort verworfen.

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

```
QoS is enabled
```

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

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

Der CoS-Wert muss der Warteschlange oder dem Schwellenwert für alle Warteschlangentypen zugewiesen werden. Die für einen 2q2t-Port-Typ definierte Zuordnung wird nicht auf einen 1p2q2t-Port angewendet. Außerdem wird die für 2q2t erstellte Zuordnung auf alle Ports angewendet, die über einen 2q2t-Warteschlangenmechanismus verfügen. Geben Sie den folgenden Befehl ein:

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

**Hinweis:** Warteschlangen werden immer nummeriert, um mit der Warteschlange mit der niedrigsten Priorität zu beginnen und mit der Warteschlange mit der höchsten Priorität, die verfügbar ist, zu enden. Hier ein Beispiel:

- Warteschlange 1 ist die WRR-Warteschlange mit niedriger Priorität.
- Warteschlange 2 ist die WRR-Warteschlange mit hoher Priorität.
- Warteschlange 3 ist die Warteschlange mit höchster Priorität.

Sie müssen diesen Vorgang für alle Warteschlangentypen wiederholen. Andernfalls behalten Sie die standardmäßige CoS-Zuweisung bei. Hier ein Beispiel für 1p2q2t:

## Konfiguration

```
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
```

## Konsolenausgabe

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

QoS tx priority queue and threshold mapped to cos successfully

Sie müssen das WRR-Gewicht für die beiden WRR-Warteschlangen konfigurieren. Geben Sie den folgenden Befehl ein:

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

*Weight\_1* bezieht sich auf die Warteschlange 1, die die WRR-Warteschlange mit niedriger Priorität sein sollte. *Weight\_1* muss immer niedriger sein als *weight\_2*. Das Gewicht kann zwischen 1 und 255 liegen. Sie können den Prozentsatz den folgenden Formeln zuweisen:

- Warteschlange 1:

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

- Warteschlange 2:

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

Sie müssen auch das Gewicht für die verschiedenen Warteschlangentypen definieren. Das Gewicht muss nicht gleich sein. Beispiel: Bei 2q2t, bei dem Warteschlange 1 30 Prozent der Zeit serviert wird und Warteschlange 2 70 Prozent der Zeit bedient wird, können Sie diesen Befehl ausführen, um das Gewicht zu definieren:

```
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.*

### Konsolenausgabe

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

```
QoS wrr ratio is set successfully
```

Sie müssen auch das Verhältnis der Übertragungswarteschlange definieren, das sich auf die Aufteilung der Puffer auf die verschiedenen Warteschlangen bezieht. Geben Sie den folgenden Befehl ein:

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

**Hinweis:** Wenn Sie drei Warteschlangen (1p2q2t) haben, müssen Sie aus Hardwaregründen die WRR-Warteschlange mit hoher Priorität und die Warteschlange mit höchster Priorität auf derselben Ebene festlegen.

### Konfiguration

```
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.*

### Konsolenausgabe

```
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
```

Wie diese Konsolenausgabe veranschaulicht, muss die Summe der Warteschlangenwerte 100 betragen. Lassen Sie den größten Teil der Puffer für die WRR-Warteschlange mit niedriger Priorität, da diese Warteschlange die meisten Puffer benötigt. Die anderen Warteschlangen

werden mit höherer Priorität bedient.

Der letzte Schritt besteht in der Konfiguration der Schwellenwerte für die WRED-Warteschlange oder die Tail-Drop-Warteschlange. Geben Sie folgende Befehle ein:

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

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

## Konfiguration

```
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*

## Konsolenausgabe

```
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
```

Der Befehl **set qos wred 1p2q2t tx queue 2 40 80** arbeitet mit dem CoS für die Grenzwertzuordnung zusammen. Wenn Sie beispielsweise die Befehle in der folgenden Liste ausgeben, stellen Sie sicher, dass Pakete mit CoS 0, 1, 2 und 3 am 1p2q2t-Port in der Übertragungsrichtung in die erste Warteschlange (die niedrige WRR-Warteschlange) gesendet werden. Wenn die Puffer in dieser Warteschlange zu 50 % gefüllt sind, beginnt WRED, Pakete mit CoS 0 und 1 zu verwerfen. Pakete mit CoS 2 und 3 werden nur verworfen, wenn die Puffer in der Warteschlange zu 90 % gefüllt sind.

```
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 lp2q2t tx 1 2 cos 3
```

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

## Überwachen der Ausgabeplanung und Überprüfen der Konfiguration

Ein einfacher Befehl zum Überprüfen der aktuellen Laufzeitkonfiguration für die Ausgabeplanung eines Ports ist **show qos info runtime mod/port**. Der Befehl zeigt folgende Informationen an:

- Die Art der Warteschlange am Port
- Die Zuordnung von CoS zu den verschiedenen Warteschlangen und Schwellenwerten
- Puffer-Freigabe
- Das WRR-Gewicht

In diesem Beispiel beträgt der WRR für die Warteschlange 1 20 % und für die Warteschlange 2 80 %.

```
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 lp1q4t (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)

Beachten Sie im nächsten Beispiel, dass die WRR-Gewichtungen nicht der Standardwert 1 sind. Die Gewichtungen wurden für die Warteschlange 1 auf die Werte 20 und für die Warteschlange 2 auf 80 festgelegt. In diesem Beispiel wird ein Datenverkehrsgenerator verwendet, um 2 GB Datenverkehr an einen Catalyst 6000 zu senden. Diese 2 GB Datenverkehr sollten den Port 1/1 passieren. Da Port 1/1 überbelegt ist, werden viele Pakete verworfen (1 Gbit/s). Der Befehl **show mac** zeigt, dass es eine Menge Ausgabeabfall gibt:

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

Betrachten Sie die verworfenen Pakete. So wird das vorgeschlagene Datenverkehrsmuster aufgeteilt:

- 1 GB Datenverkehr mit IP-Priorität 0
- 250 MB Datenverkehr mit IP-Priorität 4
- 250 MB Datenverkehr mit IP-Rangfolge 5
- 250 MB Datenverkehr mit IP-Priorität 6
- 250 MB Datenverkehr mit IP-Vorrang 7

Laut CoS-Zuordnung wird dieser Datenverkehr gesendet:

- 1 GB Datenverkehr zur Warteschlange 1 Schwellenwert 1
- 0 MB Datenverkehr zur Warteschlange 1 Schwellenwert 2
- 500 MB Datenverkehr an Warteschlange 2 Schwellenwert 1
- 250 MB Datenverkehr an Warteschleife 2 Schwellenwert 2
- 250 MB Datenverkehr an Warteschlange 3 (Warteschlange mit strikter Priorität)

Der Switch muss dem empfangenen Datenverkehr vertrauen, damit die eingehende IP-Priorität im Switch erhalten bleibt und für die Zuordnung zum CoS-Wert für die Ausgabeplanung verwendet wird.

**Hinweis:** Die standardmäßige IP-Priorität für die CoS-Zuordnung ist IP-Rangfolge gleich CoS.

Geben Sie den Befehl **show qos stat 1/1** ein, um die verworfenen Pakete und den ungefähren Prozentsatz anzuzeigen:

- An diesem Punkt werden keine Pakete in der Warteschlange 3 (CoS 5) verworfen.
- 91,85 Prozent der verworfenen Pakete sind CoS 0-Pakete in der Warteschlange 1.
- 8 Prozent der verworfenen Pakete sind CoS 4 und 6 in der Warteschlange 2, Schwellenwert 1.
- 0,15 Prozent der verworfenen Pakete sind CoS 7 in der Warteschlange 2, Schwellenwert 2.

Diese Ausgabe veranschaulicht die Verwendung des Befehls:

```
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
```

Wenn Sie das WRR-Gewicht wieder auf den Standardwert zurücksetzen, nachdem die Zähler gelöscht wurden, treten nur 1 Prozent der verworfenen Pakete in der Warteschlange 2 statt der zuvor angezeigten 8 Prozent auf:

**Hinweis:** Der Standardwert ist 5 für die Warteschlange 1 und 255 für die Warteschlange 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 : 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
```

## [Verwenden der Ausgabeplanung zur Verringerung von Verzögerungen und Jitter](#)

Das Beispiel im Abschnitt [Überwachen der Ausgabeplanung und Überprüfen der Konfiguration](#) veranschaulicht die Vorteile der Implementierung der Ausgabenplanung, wodurch bei einer Überbelegung des Ausgabeports ein Rückgang des VoIP- oder geschäftskritischen Datenverkehrs vermieden wird. Eine Überbelegung tritt selten in einem normalen Netzwerk auf, insbesondere bei einer Gigabit-Verbindung. In der Regel erfolgt die Überbelegung nur zu Spitzenzeiten des Datenverkehrs oder bei Datenverkehrsspitzen innerhalb kürzester Zeit.

Auch ohne Überbelegung kann die Ausgabeplanung in einem Netzwerk, in dem QoS durchgängig

implementiert wird, von großem Nutzen sein. Die Ausgabeplanung reduziert Verzögerungen und Jitter. In diesem Abschnitt finden Sie Beispiele dafür, wie die Ausgabeplanung Verzögerungen und Jitter reduzieren kann.

## Verzögerung reduzieren

Die Verzögerung eines Pakets wird um die Zeit erhöht, die im Puffer jedes Switches während der Wartezeit auf die Übertragung "verloren" ist. Beispielsweise wird ein kleines Sprachpaket mit einer CoS von 5 während einer großen Datensicherung oder Dateiübertragung von einem Port gesendet. Wenn Sie für den Ausgangsport keine QoS haben und davon ausgehen, dass das kleine Sprachpaket nach 10 großen 1500-Byte-Paketen in die Warteschlange gestellt wird, können Sie die für die Übertragung der 10 großen Pakete erforderliche Gigabit-Geschwindigkeit leicht berechnen:

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

Wenn dieses Paket acht oder neun Switches durchlaufen muss, während es das Netzwerk durchläuft, kann es zu einer Verzögerung von etwa 1 ms kommen. Dieser Betrag zählt nur Verzögerungen in der Ausgabewarteschlange des Switches, der im Netzwerk überschritten wird.

**Hinweis:** Wenn Sie dieselben zehn großen Pakete auf einer 10-Mbit/s-Schnittstelle in eine Warteschlange stellen müssen (z. B. bei einem IP-Telefon und einem angeschlossenen PC), wird folgende Verzögerung eingeführt:

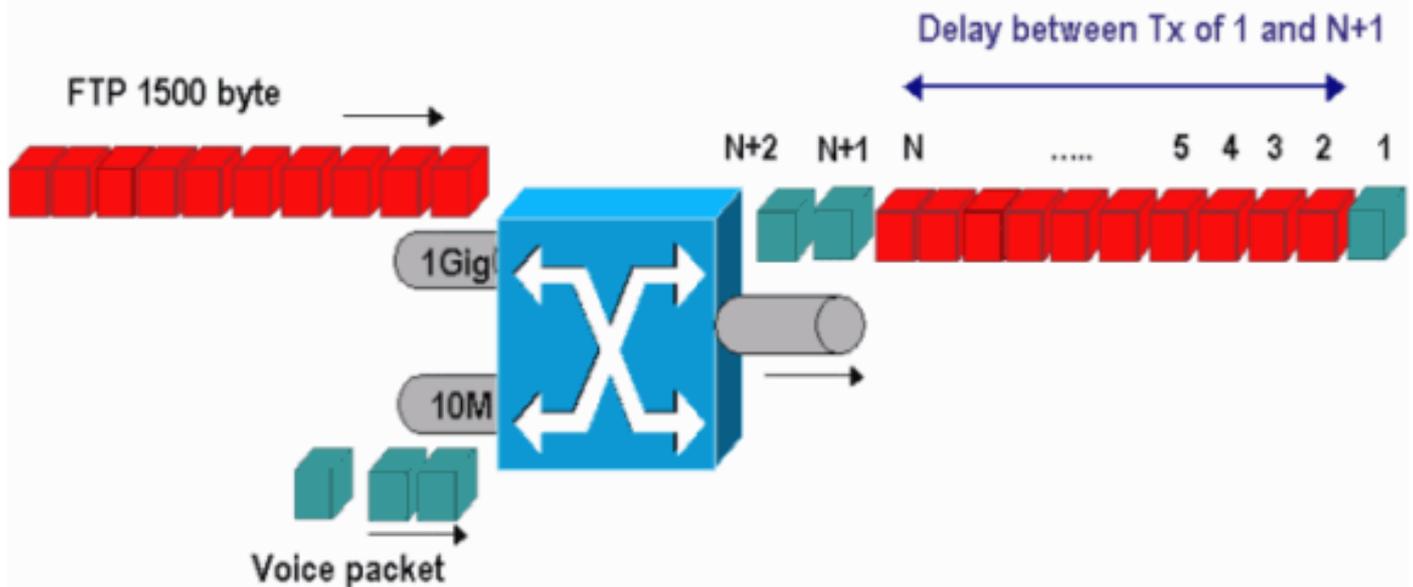
`(10 × 1500 × 8) = 120,000 bits that are transmitted in 12 ms`

Die Implementierung der Ausgabeplanung stellt sicher, dass Sprachpakete mit einer CoS von 5 in die Warteschlange mit strikter Priorität gestellt werden. Diese Platzierungen stellen sicher, dass diese Pakete vor Paketen mit einer CoS von weniger als 5 gesendet werden, wodurch die Verzögerungen verringert werden.

## Jitter reduzieren

Ein weiterer wichtiger Vorteil der Implementierung der Ausgabeplanung ist die Reduzierung von Jitter. Jitter ist die Verzögerung, die bei Paketen im selben Fluss beobachtet wird. Das Diagramm in [Abbildung 8](#) zeigt ein Beispielszenario, wie die Ausgabeplanung Jitter reduzieren kann:

### **Abbildung 8**



In diesem Szenario muss ein einzelner Ausgangsport zwei Streams senden:

- Ein Sprach-Stream, der über einen 10-Mbit/s-Ethernet-Port eingeht
- Ein FTP-Stream, der über einen 1-Gbit/s-Ethernet-Uplink eingeht

Beide Streams verlassen den Switch über denselben Ausgangsport. Dieses Beispiel zeigt, was ohne die Verwendung der Ausgabeplanung geschehen kann. Alle großen Datenpakete können zwischen zwei Sprachpaketen verschachtelt werden, wodurch beim Empfang des Sprachpakets vom gleichen Stream Jitter entsteht. Es besteht eine größere Verzögerung zwischen dem Empfang des Pakets  $n$  und des Pakets  $n+1$ , wenn der Switch das große Datenpaket überträgt. Die Verzögerung zwischen  $n+1$  und  $n+2$  ist jedoch vernachlässigbar. Dies führt zu Jitter im Sprachdatenverkehrsstrom. Sie können dieses Problem mit einer Warteschlange mit strikter Priorität leicht vermeiden. Stellen Sie sicher, dass der CoS-Wert der Sprachpakete der Warteschlange mit strikter Priorität zugeordnet ist.

## Zugehörige Informationen

- [QoS-Ausgabeplanung für Catalyst Switches der Serien 6500/6000 mit Cisco IOS-Systemsoftware](#)
- [Quality of Service auf Catalyst Switches der Serie 6000](#)
- [Support-Seiten für LAN-Produkte](#)
- [Support-Seite für LAN-Switching](#)
- [Technischer Support und Dokumentation - Cisco Systems](#)