



Cisco TelePresence Network Systems 2.0 Design Guide

Cisco Validated Design

August 1, 2008

Cisco Validated Designs for deploying point-to-point and multipoint Cisco TelePresence 500, 1000, 3000, and 3200 systems in enterprise campus and branch, WAN, and VPN networks.

Americas Headquarters

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

Customer Order Number: OL-14133-01

Cisco Validated Design

The Cisco Validated Design Program consists of systems and solutions designed, tested, and documented to facilitate faster, more reliable, and more predictable customer deployments. For more information visit www.cisco.com/go/validateddesigns.

ALL DESIGNS, SPECIFICATIONS, STATEMENTS, INFORMATION, AND RECOMMENDATIONS (COLLECTIVELY, "DESIGNS") IN THIS MANUAL ARE PRESENTED "AS IS," WITH ALL FAULTS. CISCO AND ITS SUPPLIERS DISCLAIM ALL WARRANTIES, INCLUDING, WITHOUT LIMITATION, THE WARRANTY OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT OR ARISING FROM A COURSE OF DEALING, USAGE, OR TRADE PRACTICE. IN NO EVENT SHALL CISCO OR ITS SUPPLIERS BE LIABLE FOR ANY INDIRECT, SPECIAL, CONSEQUENTIAL, OR INCIDENTAL DAMAGES, INCLUDING, WITHOUT LIMITATION, LOST PROFITS OR LOSS OR DAMAGE TO DATA ARISING OUT OF THE USE OR INABILITY TO USE THE DESIGNS, EVEN IF CISCO OR ITS SUPPLIERS HAVE BEEN ADVISED OF THE POSSIBILITY OF SUCH DAMAGES.

THE DESIGNS ARE SUBJECT TO CHANGE WITHOUT NOTICE. USERS ARE SOLELY RESPONSIBLE FOR THEIR APPLICATION OF THE DESIGNS. THE DESIGNS DO NOT CONSTITUTE THE TECHNICAL OR OTHER PROFESSIONAL ADVICE OF CISCO, ITS SUPPLIERS OR PARTNERS. USERS SHOULD CONSULT THEIR OWN TECHNICAL ADVISORS BEFORE IMPLEMENTING THE DESIGNS. RESULTS MAY VARY DEPENDING ON FACTORS NOT TESTED BY CISCO.

CCDE, CCENT, Cisco Eos, Cisco Lumin, Cisco Nexus, Cisco StadiumVision, Cisco TelePresence, the Cisco logo, DCE, and Welcome to the Human Network are trademarks; Changing the Way We Work, Live, Play, and Learn and Cisco Store are service marks; and Access Registrar, Aironet, AsyncOS, Bringing the Meeting To You, Catalyst, CCDA, CCDP, CCIE, CCIP, CCNA, CCNP, CCSP, CCVP, Cisco, the Cisco Certified Internetwork Expert logo, Cisco IOS, Cisco Press, Cisco Systems, Cisco Systems Capital, the Cisco Systems logo, Cisco Unity, Collaboration Without Limitation, EtherFast, EtherSwitch, Event Center, Fast Step, Follow Me Browsing, FormShare, GigaDrive, HomeLink, Internet Quotient, IOS, iPhone, iQ Expertise, the iQ logo, iQ Net Readiness Scorecard, iQuick Study, IronPort, the IronPort logo, LightStream, Linksys, MediaTone, MeetingPlace, MeetingPlace Chime Sound, MGX, Networkers, Networking Academy, Network Registrar, PCNow, PIX, PowerPanels, ProConnect, ScriptShare, SenderBase, SMARTnet, Spectrum Expert, StackWise, The Fastest Way to Increase Your Internet Quotient, TransPath, WebEx, and the WebEx logo are registered trademarks of Cisco Systems, Inc. and/or its affiliates in the United States and certain other countries.

All other trademarks mentioned in this document or Website are the property of their respective owners. The use of the word partner does not imply a partnership relationship between Cisco and any other company. (0807R)

Cisco TelePresence Network Systems 2.0 Design Guide
Copyright © 2008 Cisco Systems, Inc. All rights reserved.



CONTENTS

CHAPTER 1

Cisco TelePresence Solution Overview	1-1
Cisco TelePresence System 3000	1-1
Cisco TelePresence System 3200	1-2
Cisco TelePresence System 1000	1-3
Cisco TelePresence System 500	1-4
Cisco TelePresence Codecs	1-5
Industry-Leading Audio and Video Support	1-7
Video Resolutions and Compression Formats	1-7
Resolution	1-7
Frame Rate	1-8
Compression	1-8
Audio Resolution and Compression Formats	1-8
Frequency Spectrum	1-9
Spatiality	1-9
Compression	1-9
Cisco TelePresence Manager	1-9
Cisco Unified 7975G IP Phone	1-11
Cisco TelePresence Multipoint Solutions	1-12

CHAPTER 2

Connecting the Endpoints	2-1
Overview	2-1
Connecting a CTS-1000 System	2-1
Connecting a CTS-500 System	2-3
Connecting a CTS-3000 System	2-3
Connecting a CTS-3200 System	2-5
Cisco TelePresence Network Interaction	2-8

CHAPTER 3

TelePresence Network Deployment Models	3-1
Introduction	3-1
Intra-Campus Deployment Model	3-1
Intra-Enterprise Deployment Model	3-2

- Cisco Powered Networks 3-3
- Point-to-Point versus Multipoint 3-3
- Inter-Enterprise/Business-to-Business Deployment Model 3-4
- Hosting and Management Options 3-5
- TelePresence Phases of Deployment 3-5

CHAPTER 4

Quality of Service Design for TelePresence 4-1

- Overview 4-1
- Defining the Strategic Business Objective for QoS for TelePresence 4-2
- Analyzing the Service Level Requirements of TelePresence 4-3
 - TelePresence Bandwidth Requirements 4-3
 - Burst Requirements 4-6
 - TelePresence Latency Requirements 4-7
 - TelePresence Jitter Requirements 4-9
 - TelePresence Loss Requirements 4-10
- Tactical QoS Design Best Practices for TelePresence 4-12
 - Relevant Industry Standards and Recommendations 4-13
 - RFC 2474 Class Selector Code Points 4-13
 - RFC 2597 Assured Forwarding Per-Hop Behavior Group 4-13
 - RFC 3246 An Expedited Forwarding Per-Hop Behavior 4-13
 - RFC 3662 A Lower Effort Per-Domain Behavior for Differentiated Services 4-13
 - Cisco's QoS Baseline 4-13
 - RFC 4594 Configuration Guidelines for DiffServ Classes 4-14
 - Classifying TelePresence 4-17
 - Policing TelePresence 4-18
 - Queuing TelePresence 4-19
 - Shaping TelePresence? 4-20
 - Compressed RTP (cRTP) with TelePresence 4-20
 - Link Fragmentation and Interleaving (LFI) with TelePresence 4-21
 - GRE/IPSec Tunnels with TelePresence 4-21
- Place in the Network TelePresence QoS Design 4-21

CHAPTER 5

Campus QoS Design for TelePresence 5-1

- Overview 5-1
- Access Edge Switch Port QoS Considerations 5-1
- Campus Inter-Switch Link QoS Considerations 5-5
- TelePresence Campus QoS Designs 5-6
 - Catalyst 3560G/3750G and 3650-E/3750E 5-7

Catalyst 4500 and 4948	5-13
Catalyst 6500	5-17
Ingress Queuing Design—1Q2T	5-21
Egress Queuing Design—1P2Q2T	5-22
Egress Queuing Design—1P3Q8T	5-23
Egress Queuing Design—1P7Q8T	5-25
Egress Queuing Design—1P7Q4T (DSCP-to-Queue)	5-28

CHAPTER 6**Branch QoS Design for TelePresence 6-1**

TelePresence Branch QoS Design Overview	6-1
LLQ versus CBWFQ Considerations	6-1
Campus WAN/VPN Block Considerations	6-7
TelePresence Branch LAN Edge	6-8
TelePresence Branch LAN Edge QoS Design Considerations	6-8
TelePresence Branch LAN Edge QoS Designs	6-11
TelePresence Branch WAN Edge	6-11
TelePresence Branch WAN Edge Design Considerations	6-11
TelePresence Branch WAN Edge QoS Design	6-11
TelePresence Branch WAN Edge LLQ Policy	6-11
TelePresence Branch WAN Edge CBWFQ Policy	6-14
TelePresence Branch T3/DS3 WAN Edge Design	6-14
TelePresence Branch OC3-POS WAN Edge Design	6-18
TelePresence Branch IPsec VPN Edge	6-22
TelePresence Branch IPsec VPN Edge Considerations	6-22
TelePresence Branch IPsec VPN Edge QoS Design	6-24
TelePresence Branch MPLS VPN	6-26
TelePresence Branch MPLS VPN Edge Considerations	6-27
TelePresence Branch MPLS VPN QoS Designs	6-32
TelePresence 4-Class MPLS VPN SP Model QoS Design	6-32
TelePresence 6-Class MPLS VPN SP Model QoS Design	6-37
TelePresence Sub-Line Rate Ethernet Access QoS Designs	6-39

CHAPTER 7**Call Processing Overview 7-1**

Overview	7-1
Call Processing Components	7-1
TelePresence Endpoint Interface to CUCM (Line-Side SIP)	7-3
TelePresence Multipoint Switch Interface to CUCM (Trunk-Side SIP)	7-3
TelePresence Endpoint Device Registration	7-4

Call Setup 7-4
 Call Teardown 7-7
 Firewall and NAT Considerations 7-8

CHAPTER 8

Capacity Planning and Call Admission Control 8-1

Overview 8-1
 Manual Capacity Planning 8-1

CHAPTER 9

Call Processing Deployment Models 9-1

Overview 9-1
 Dial-Plan Recommendations 9-1
 Single-Site Call Processing Model 9-2
 Call Admission Control 9-3
 Multi-Site WAN with Centralized Call Processing Model 9-4
 Call Admission Control 9-4

CHAPTER 10

Cisco TelePresence Multipoint Solution Essentials 10-1

Overview of Multipoint Conference Technologies 10-1
 Components of the Cisco TelePresence Multipoint Solution 10-3
 Overview 10-3
 Multipoint Virtual Meetings Which Include Only CTS Endpoints 10-4
 Multipoint Virtual Meetings Which Also Include Traditional Video Conferencing Systems 10-5
 Cisco TelePresence Multipoint Switch Overview 10-5
 Meeting Types 10-6
 Static Meetings 10-6
 Ad Hoc Meetings 10-8
 Scheduled Meetings 10-8
 CTMS Meeting Features 10-9
 Switching Policy 10-10
 Maximum Number of Rooms 10-13
 Video Announce 10-13
 Lock Meeting 10-14
 Quality 10-14
 VIP Mode 10-14
 Multipoint Resources 10-14
 Geographical Resource Management 10-16
 Quality of Service 10-17
 Meeting Security 10-18

Administrative Access Control	10-18
Meeting Access Control	10-19
Meeting Confidentiality	10-20
Meeting Management	10-20

CHAPTER 11**Cisco Multipoint Technology and Design Details 11-1**

Audio and Video Flows In A Multipoint TelePresence Design	11-1
Flow Control Overview	11-1
Audio and Video Positions	11-1
Audio to the CTMS in a Multipoint TelePresence Meeting	11-3
Calculating the Amount of Audio Traffic to the CTMS	11-4
Audio From the CTMS in a Multipoint Meeting	11-5
Calculating the Amount of Audio Traffic from the CTMS	11-5
Video in a Multipoint TelePresence Meeting	11-6
Camera Video Input	11-7
Auxiliary Video Input	11-7
Calculating the Amount of Video Traffic to the CTMS	11-7
Calculating the Amount of Video Traffic From the CTMS	11-9
Total Traffic to and from the CTMS	11-10
Video Switchover Delay	11-10
Overview of TelePresence Video on the Network	11-11
Deployment Models	11-15
Centralized Deployment	11-15
Deployment Considerations	11-16
Distributed Deployment	11-19
Deployment Considerations	11-20
Positioning of the CTMS within the Campus or Branch	11-23
Network Requirements	11-25
Latency	11-26
Bandwidth	11-27
Estimating Burst Sizes within Multipoint TelePresence Calls	11-32
Causes of Bursts within Multipoint TelePresence Calls	11-32
Bursts Due to I-Frame Replication	11-33
Location of the CTS Endpoints	11-34
Type of CTS Endpoints	11-35
Calculating Burst Sizes Due to I-Frame Replication	11-37
Other Considerations	11-39
Bursts due to the Auxiliary Video Input	11-40
Burst Estimation Due to Auxiliary Video Replication	11-41

- Other Considerations 11-42
- Normal P-Frame Video 11-42
 - Location of the CTS Endpoints 11-43
 - Burst Estimation Due to P-Frame Replication 11-43
 - Other Considerations 11-44

CHAPTER 12

Cisco TelePresence Multipoint Solution Circuit and Platform Recommendations 12-1

- Recommendations for Multipoint over WAN Circuits 12-1
 - Cisco Router and Switch Platforms Tested 12-1
 - TelePresence over Dedicated WAN Circuits 12-2
 - Dedicated T3, E3, and OC-3 POS Circuits 12-2
 - Multipoint Over Dedicated OC-12 and OC-48 POS Circuits 12-7
 - Recommendations for Multipoint Over Converged WAN Circuits 12-7
 - Multipoint TelePresence WAN Edge LLQ Policy 12-8
 - Multipoint TelePresence Branch WAN Edge CBWFQ Policy 12-12
- Multipoint TelePresence over MPLS Circuits with Ethernet Handoff 12-14
 - Cisco Router and Switch Platforms Tested 12-14
 - Multipoint TelePresence over Dedicated MPLS Circuits 12-15
 - Multipoint TelePresence over Converged MPLS Circuits 12-18
- Platform and Linecard Test Results and Recommendations 12-21
 - Multipoint TelePresence over Switch Linecards and Platforms within the LAN 12-21
 - Multipoint TelePresence over WAN Switch and Router Platforms 12-24

CHAPTER 13

Internal Firewall Deployments with Cisco TelePresence 13-1

- Overview 13-1
- Cisco Firewall Platforms 13-3
- Firewall Deployment Options 13-3
 - Transparent versus Routed Mode 13-3
 - Equal versus Unequal Interface Security Levels 13-4
 - Network Address Translation (NAT) 13-5
 - Application Layer Protocol Inspection 13-7
- TelePresence Protocol Requirements 13-7
 - Device Provisioning Flows 13-7
 - Dynamic Host Configuration Protocol (DHCP) 13-7
 - Domain Name System (DNS) 13-9
 - Configuration Download Protocols 13-10
 - Call Scheduling and Services Flows 13-12
 - Call Signaling Flows 13-13
 - Media Flows 13-13

Point-to-Point TelePresence Calls	13-13
Multipoint TelePresence Calls	13-15
Management Flows	13-16
HTTP, HTTPS, and SSH	13-16
SNMP Traps	13-16
ESE Firewall Test Results	13-17
Example Firewall Configuration	13-18

CHAPTER 14**Cisco Services for Cisco TelePresence 14-1**

Challenge	14-2
Solution	14-2
The Cisco TelePresence Planning, Design, and Implementation Service	14-2
The Cisco TelePresence Essential Operate Service	14-4
The Cisco TelePresence Select Operate Service	14-4
The Cisco TelePresence Remote Assistance Service	14-5
Benefits	14-5
Why Cisco Services	14-6
For More Information	14-6



CHAPTER 1

Cisco TelePresence Solution Overview

The Cisco TelePresence suite of virtual meeting solutions consists of the products and capabilities described in the following sections.

Cisco TelePresence System 3000

The Cisco TelePresence System 3000 (CTS-3000) is designed for business meetings, with up to six participants per room. It consists of:

- Three 65" high definition plasma displays
- Three high definition cameras
- Three wide band microphones and speakers
- A lighting shroud integrated around a purpose built meeting room table

Customers must furnish their own chairs. A Cisco 7975G IP phone is used to launch, control, and end the meeting.

Figure 1-1 Cisco TelePresence System 3000

221627

Participants are displayed life size with two participants per screen/table segment and multi-channel, discrete, full-duplex audio with echo cancellation per channel that appears to emanate from the person speaking. The unique table design also provides power and Ethernet ports in each table leg, so users do not have to hunt for power and network connections during the meeting. A projector is integrated under the middle section of the table for convenient viewing of PC graphics on the panel below the plasma displays. An optional WolfVision® document camera (not shown) may be installed in the ceiling so that objects and documents placed on the table surface may be viewed as well.

The CTS-3000 is represented by the icon in [Figure 1-2](#).

Figure 1-2 CTS-3000 Icon

Cisco TelePresence System 3200

The Cisco TelePresence System 3200 (CTS-3200) is designed for large group and cross-functional team meetings of up to 18 participants per room. It consists of:

- Three 65" high definition plasma displays
- Three high definition cameras
- Nine wide band microphones and three speakers
- A lighting shroud integrated around a purpose built meeting room table

Customers must furnish their own chairs. A Cisco 7975G IP phone is used to launch, control, and end the meeting.

Figure 1-3 Cisco TelePresence System 3200



The CTS-3200 extends the CTS-3000 by adding a second row of seating which includes a purpose built meeting room table and six additional wide band microphones. The CTS-3200 also provides an option for supporting up to three additional additional graphics displays for convenient viewing of PC graphics via an HDMI splitter.

The CTS-3200 is represented by the same icon as the CTS-3000 shown in [Figure 1-2](#).

Cisco TelePresence System 1000

The Cisco TelePresence System 1000 (CTS-1000) is designed for small group meetings and one-on-one conversations, seating up to two participants per room. It consists of:

- One 65" high definition plasma display
- One high definition camera
- One wide band microphone and speaker
- A lighting shroud integrated over the display

The customer must furnish their own meeting room table and chairs. A Cisco IP phone is used to launch, control, and end the meeting.

Figure 1-4 Cisco TelePresence System 1000



Participants are displayed life size with two participants per screen/table segment and full-duplex audio with echo cancellation that appears to emanate from the person speaking. An optional NEC® LCD display (not shown) may be installed on the table or on the wall for convenient viewing of PC graphics. An optional WolfVision® document camera (not shown) may be installed on the table so that objects and documents placed on the table surface may be viewed as well.

The CTS-1000 is represented by the icon in [Figure 1-5](#).

Figure 1-5 CTS-1000 Icon



Cisco TelePresence System 500

The Cisco TelePresence System 500 (CTS-500) is designed with a smaller form factor and streamlined footprint to fit easily into private offices or public locations. It consists of an integrated 37” display, camera, microphone, speakers, and lighting suitable for private offices.

Figure 1-6 Cisco TelePresence System 500



The CTS 500 delivers the same superior video and audio quality as the rest of the Cisco TelePresence portfolio: 1080p video, wide-band audio, one-button-to-push call initiation, multipoint capabilities, scheduling through existing groupware applications, ad hoc calling support, and interoperability with traditional, standards-based videoconferencing systems. It is available in three configurations:

- Free-standing pedestal
- Wall mount
- Table top

The CTS-500 is represented by the icon in [Figure 1-7](#).

Figure 1-7 CTS-500 Icon



Cisco TelePresence Codecs

One of the goals of Cisco TelePresence is to hide the technology from the user so that participants experience the meeting, not the technology. Hidden underneath the plasma displays of the CTS-3200, CTS-3000, and CTS-1000 solutions are the Cisco TelePresence codecs. The CTS-3000 and CTS-3200 consist of one primary codec and two secondary codecs. The CTS-1000 and CTS-500 consist of a single primary codec. An additional optional codec is available as an upgrade for high-speed (30 frames per second) auxiliary video input.

Figure 1-8 Cisco TelePresence Codec



The codec is the engine which drives the entire Cisco TelePresence solution. All displays, cameras, microphones, and speakers connect to it and it communicates with the network and handles all audio and video processing. The codec runs a highly-integrated version of the Linux operating system on an embedded Compact Flash module and is managed via Secure Shell (SSH) and Hyper-Text Transfer Protocol over Secure Sockets Layer (HTTPs). These codecs make the Cisco TelePresence solutions an integrated part of Cisco Unified Communications by leveraging established techniques for network automation and Quality of Service (QoS), such as:

- Cisco Discovery Protocol (CDP) and 802.1Q for discovery and assignment to the appropriate Virtual LAN (VLAN).
- 802.1p and Differentiated Services Code Point (DSCP) for QoS.
- Automated provisioning of configuration and firmware from Cisco Unified Communications Manager.
- Session Initiation Protocol (SIP) for all call signaling communications.

From an administrator's perspective, the entire Cisco TelePresence virtual meeting room appears as a single SIP endpoint on Cisco Unified Communications Manager. It is managed using tools and methodologies that are similar to those used for Cisco Unified IP Phones.

The Cisco TelePresence Codec is represented by the icon in [Figure 1-9](#).

Figure 1-9 Cisco TelePresence Codec Icon



Industry-Leading Audio and Video Support

Cisco TelePresence utilizes industry-leading 1080p high-definition video resolution and 48kHz wide-band spatial audio. 720p high-definition is also supported for sites with restricted bandwidth availability.

Video Resolutions and Compression Formats

The Cisco TelePresence displays and cameras natively support 1080p resolution and utilize digital media interfaces to connect to the Cisco TelePresence codecs. This ensures the integrity of the video signal from end to end by eliminating the need for any digital/analog conversion.

Inside the Cisco TelePresence codecs an onboard array of Digital Signal Processors (DSPs) encode the digital video signal from the cameras into Real-Time Transport Protocol (RTP) packets using the H.264 encoding and compression standard. The Cisco TelePresence codecs can encode the video from the cameras at 1080p or 720p.

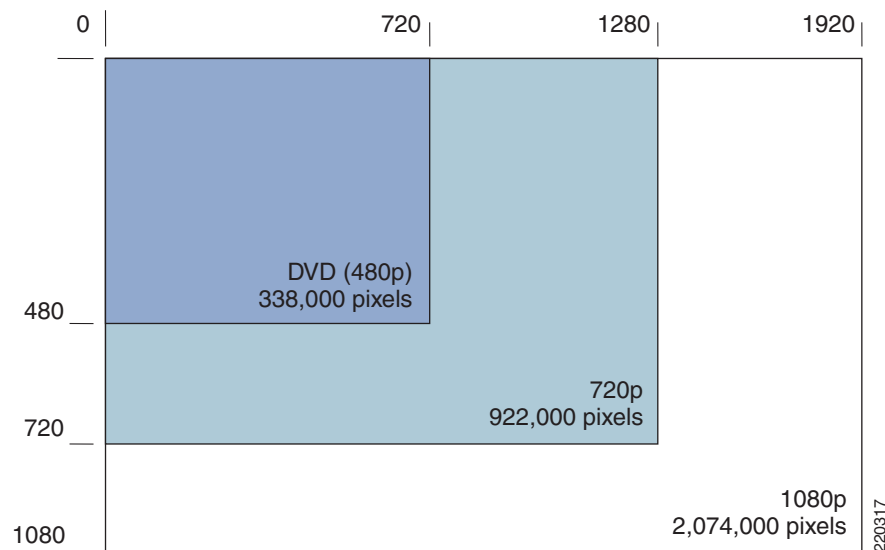
The quality of the video enjoyed by the meeting participants is a function of three variables:

- Resolution (i.e., number of pixels within the image)
- Frame rate (how often those pixels are re-drawn on the display)
- Degree of compression applied to the original video signal

Resolution

1080p provides the highest quality video image currently available on the market, supplying a resolution of 1920 x 1080 and 2,074,000 pixels per frame. 720p provides a resolution of 1280 x 720 and 922,000 pixels per frame. Compared with today's DVD standard video (480p) with a resolution of 720 x 480 and 338,000 pixels per frame, you can see the dramatic increase in resolution and pixel count. [Figure 1-10](#) illustrates the difference between these three resolutions.

Figure 1-10 Video Resolutions



Frame Rate

The frame rate of the displayed video directly corresponds to how motion within the video is perceived by the participants. To maintain excellent motion handling, the Cisco TelePresence System encodes video from the cameras at 30 frames per second (30fps or 30Hz). In addition, the codec video output signal to the plasma displays utilizes progressive-scan technology in order to eliminate any unpleasant visual artifacts that sometimes result from interlaced scan technology.

Compression

Note that 1080p video uncompressed is approximately 1.5 Gbps. The Cisco TelePresence Codecs must take this native video received from the cameras and compress it to a more feasible bandwidth value in as little time as possible. As mentioned above, they achieve this by utilizing an array of DSPs to compress the original 1.5 Gbps video from each camera down to under 4 Mbps (per camera), representing a compression ratio of over 99%, and they achieve this in under 90ms. To provide maximum flexibility, the customer is provided with some amount of control over how much compression is applied. For each of the two resolution formats supported (1080p and 720p), the Cisco TelePresence System supports three quality levels. Each quality level is really a function of the degree of compression applied, and has a corresponding bandwidth value. For simplicity, these three levels are referred to as “good,” “better,” and “best.” The “best” quality level has the least amount of compression applied and therefore requires the most bandwidth, while the “good” quality level has the most amount of compression applied and requires the least amount of bandwidth.

Taking the three variables described above—resolution, frame rate, and the degree of compression applied—[Table 1-1](#) illustrates the different quality settings supported by the Cisco TelePresence System and the requisite bandwidth required for each quality setting.

Table 1-1 Resolution, Quality, and Bandwidth Settings Supported (Video Only)

Resolution	1080p			720p		
	Best	Better	Good	Best	Better	Good
Quality Level						
Frame Rate	30	30	30	30	30	30
Bandwidth Required	4Mbps	3.5Mbps	3Mbps	2.25Mbps	1.5Mbps	1Mbps

These bandwidth values apply per camera. Therefore, a CTS-3000 which has three cameras and three displays, running at 1080p resolution at the “best” quality level, requires 12Mbps of video bandwidth, whereas a CTS-1000 requires 4Mbps of video bandwidth. These bandwidth values do not include the audio channels, the auxiliary video channel for displaying PC graphics and document camera images, or network layer overhead. Therefore, a more complete bandwidth table is [Table 4-1](#).

Audio Resolution and Compression Formats

The Cisco TelePresence System utilizes advanced microphone, speaker, and audio encoding technologies to preserve the quality and directionality of the audio so that it appears to emanate from the location of the person speaking at the same volume as it would be heard if that person were actually sitting across the table from you. Specifically, wideband spatial audio and multi-channel, full-duplex sound provides excellent voice projection and helps enable multiple simultaneous conversations, just like what typically occurs during an in-person meeting. Specially designed microphones eliminate sound interference.

The quality of the audio enjoyed by the meeting participants is a function of three variables:

- Frequency spectrum and decibel levels captured by the microphones
- Spatiality (i.e., directionality) of the audio
- Degree of compression applied to the original audio signal

Frequency Spectrum

The Cisco TelePresence microphones are designed to capture a 48kHz frequency spectrum of audio in a directional pattern that focuses on the people sitting directly in front of it and are geared to the decibel levels of human speech. Filters are designed into the microphones to eliminate interference from GSM and GPRS cellular signals and to eliminate certain frequencies generated by machinery such as the fans found in laptop computers and Heating, Ventilation, and Air Conditioning (HVAC) systems. Echo cancellation technology is built into the Cisco TelePresence Codec to eliminate cross-talk and double-talk.

The Cisco TelePresence speakers are designed to reproduce the same rich frequency spectrum and decibel level of human speech.

Spatiality

To preserve the spatiality (i.e., directional perception) of the audio, the CTS-3000 employs three individual microphones placed at specific locations of the virtual table, along with three individual speakers located under each display. The CTS-3200 preserves spatiality of the audio by adding another six individual microphones placed at specific locations on the table for the second row of participants.

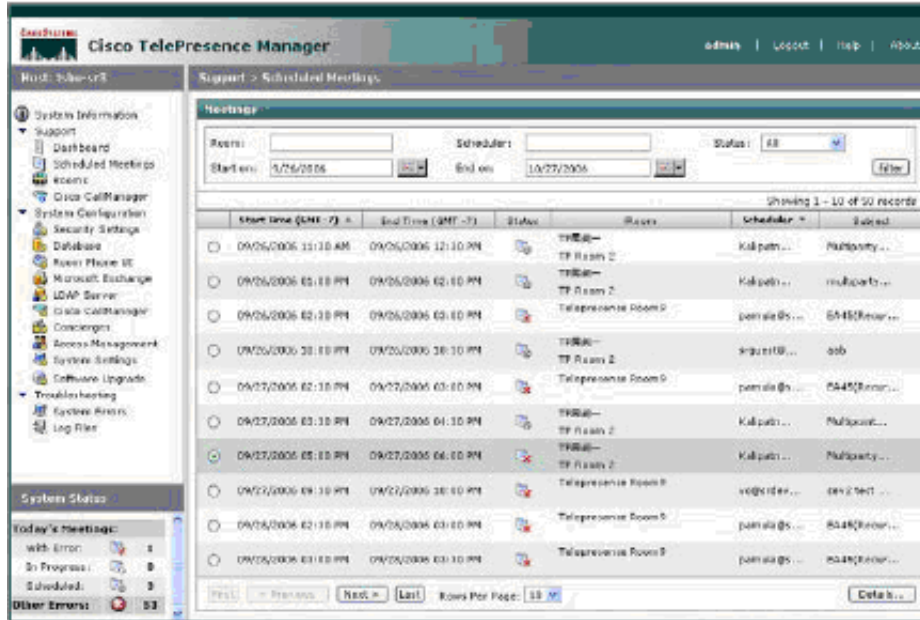
Compression

Inside the Cisco TelePresence Codecs an onboard array of DSPs encode the audio signal from the microphones into RTP packets using the Advanced Audio Coding-Low Delay (AAC-LD) encoding and compression standard. The resulting bandwidth required to transport the audio signals between the systems is 64kbps per microphone. Therefore, a CTS-3000 which has three microphones and speakers requires 192kbps of audio bandwidth, whereas the CTS-1000 requires 64kbps of audio bandwidth. Note that the Cisco TelePresence System also supports a fourth auxiliary audio channel which is used to transmit audio from a PC (used in conjunction with the projector when displaying PC graphics) or from an audio-only participant which is conferenced into the meeting using the Conference/Join softkey on the Cisco 7975G IP Phone (also known as the Audio Add-In feature). Therefore, a CTS-3000 can transmit and receive up to 256kbps of audio, as detailed in [Table 4-1](#). The CTS-1000 transmits up to 128kbps of audio, but can receive up to 256kbps when participating in a meeting with a CTS-3000 (in such a configuration, the CTS-1000 receives three separate [64 kbps] primary audio streams from the CTS-3000, as well as a potentially additional [64 kbps] auxiliary audio stream). Note that none of these bandwidth numbers include network overhead.

Cisco TelePresence Manager

Cisco TelePresence Manager (CTS-MAN) simplifies the scheduling and management of Cisco TelePresence virtual meeting room solutions. CTS-MAN is a Linux-based appliance running on a Cisco 7800 Series Media Convergence Server platform. It is the middleware glue between Cisco Unified Communications Manager, the Cisco TelePresence meeting rooms, and the customer's groupware calendaring and scheduling application (e.g., Microsoft Exchange/Outlook).

Figure 1-11 Cisco TelePresence Manager



CTS-MAN collects information about Cisco TelePresence systems from Cisco Unified Communications Manager and associates those systems to their physical location or conference room as defined in the customer's Microsoft Active Directory and Microsoft Exchange and IBM Domino.¹ This allows users to schedule Cisco TelePresence meetings using their Microsoft Outlook group calendar and have that schedule automatically sent to the Cisco TelePresence systems involved in the call. Hence users can launch the Cisco TelePresence call with the push of one button, by simply selecting their meeting from the list of meetings shown on the Cisco Unified 7975G IP phone in the meeting room.

CTS-MAN is managed via SSH, HTTPs, and SNMP. From an administrator's perspective, CTSMGR is managed using tools and methodologies that are similar to those used with a Cisco Unified Communications Manager server.

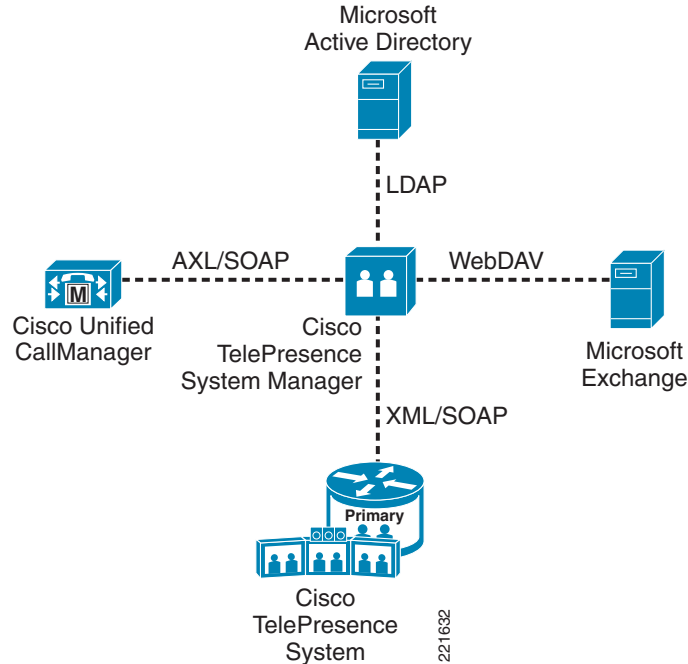
CTS-MAN communicates with Cisco Unified Communications Manager using Application XML Layer/Simple Object Access Protocol (AXL/SOAP) and Computer Telephony Integration/Quick Buffer Encoding (CTI/QBE).

CTS-MAN communicates with Microsoft Active Directory and Microsoft Exchange using Light-Weight Directory Access Protocol (LDAP) and Web-Based Distributed Authoring and Versioning (WebDAV) standards.

CTS-MAN communicates with the Cisco TelePresence Systems using eXtensible Markup Language/Simple Object Access Protocol (XML/SOAP).

1. Cisco TelePresence Manager supports Microsoft Active Directory 2003 and Microsoft Exchange 2003 and 2007. Cisco TelePresence Manager also supports IBM Domino 7.0 and Notes 6.5 and 7.0.

Figure 1-12 Cisco TelePresence Manager Connectivity



Cisco Unified 7975G IP Phone

To further enhance the meeting participants' experience of the meeting, cumbersome hand-held remote controls are eliminated, the cameras are fixed in their positions (no panning, tilting, or zooming controls), and the microphones are fixed in their positions on the table. There are virtually no moving parts or user interfaces that users must master to use a Cisco TelePresence meeting room.

Rather, the Cisco TelePresence meeting room solutions use a Cisco Unified 7975G IP phone, conveniently located on the table, to launch, control, and conclude meetings. For simplicity, the IP Phone is referred to as a 7975G IP phone. This makes Cisco TelePresence as easy to use as a telephone. Using the high-resolution touch-screen display of the Cisco Unified 7975G IP phone, the user simply dials the telephone number of the Cisco TelePresence room with which they wish to have a meeting and the call is connected. Softkey menu buttons on the phone allow the user to place the call on hold or conference in an audio-only participant. When used in conjunction with Cisco TelePresence Manager, the schedule of meetings for the day are displayed on the phone and the user simply touches the appropriate location on the screen to launch that scheduled meeting.

Figure 1-13 Cisco Unified 7975G IP Phone

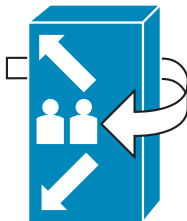


Cisco TelePresence Multipoint Solutions

To enable Cisco TelePresence meetings between more than two rooms, a Cisco TelePresence Multipoint Switch (CTMS) is required. The Cisco TelePresence Multipoint Switch is a purpose-built Linux-based appliance running on a Cisco 7800 Series Media Convergence Server platform. It provides high-capacity, low-latency multipoint switching for Cisco TelePresence only. Details about the Cisco TelePresence Multipoint Solution can be found in [Chapter 10, “Cisco TelePresence Multipoint Solution Essentials,”](#) [Chapter 11, “Cisco Multipoint Technology and Design Details,”](#) and [Chapter 12, “Cisco TelePresence Multipoint Solution Circuit and Platform Recommendations.”](#)

The CTMS is represented by the icon in [Figure 1-14](#).

Figure 1-14 CTMS Icon





CHAPTER 2

Connecting the Endpoints

Overview

As discussed in [Chapter 1, “Cisco TelePresence Solution Overview,”](#) there are many elements to Cisco TelePresence endpoint systems, including:

- TelePresence codecs (primary and secondary)
- Cisco Unified 7975G IP phone
- 65” plasma displays
- Cameras
- Microphones
- Speakers
- Auxiliary audio devices
- Auxiliary video devices

There are other elements, such as mounting brackets, furniture, cables, and power cords; the full assembly and connectivity instructions are covered in detail in the documentation.

The focus of this chapter is to provide an overview of how these main system elements are interconnected within CTS-1000, CTS-500, CTS-3000, and CTS-3200 systems, as well as how these interact with the network infrastructure. Such an overview helps lay a foundational context for the design chapters that follow.

Connecting a CTS-1000 System

The CTS-1000 includes:

- One Cisco TelePresence codec (a primary codec)
- One Cisco Unified 7975G IP phone
- One 65” plasma display
- One high-definition camera
- One microphone
- One speaker
- One input for auxiliary audio

- One input for auxiliary video which can be used for a document camera or PC

The Cisco TelePresence primary codec is the center of the CTS systems. Essentially, all components connect to it and it, in turn, connects to the network infrastructure.

Specifically, the Cisco Unified 7975G IP phone connects to the TelePresence primary codec via an RJ-45 cable that provides it network connectivity and 802.3af Power-over-Ethernet (PoE).

Another RJ-45 cable connects from the TelePresence primary codec to the camera, providing the camera with 802.3af PoE. A second cable from the primary codec to the camera provides video connectivity.

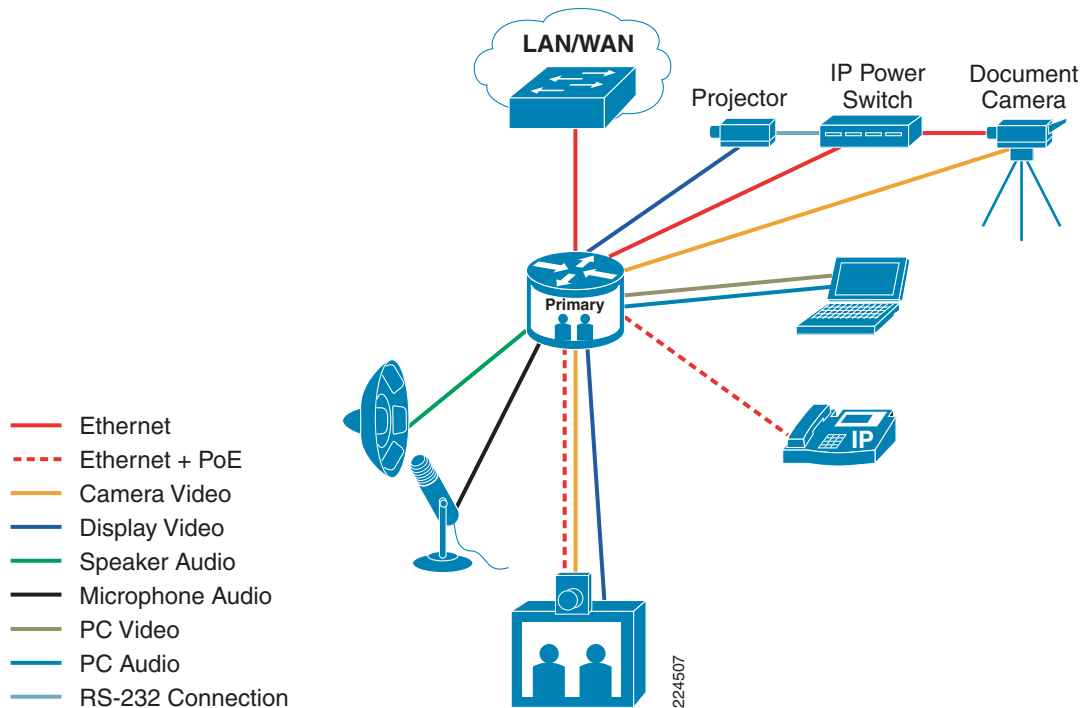
A video cable also connects the primary codec to the 65" plasma display. This cable is essentially an High Definition Multimedia Interface (HDMI) cable, but with a proprietary element for carrying management information instead of audio signals (as the audio signals are processed independently by the master codec).

Additionally, a speaker cable and a microphone cable connect the speaker and microphone to the primary codec, respectively.

The primary codec also has inputs for auxiliary audio and auxiliary video. Auxiliary video can come from a PC connection or from a document camera connection. An IP power switch (IPS) provides control for the on/off function of the document camera, attached projector, as well as the lighting shroud of the CTS unit via an Ethernet connection.

Finally, an RJ-45 cable provides 10/100/1000 Ethernet connectivity from the primary codec to the network infrastructure. These interconnections for a CTS-1000 system are illustrated in [Figure 2-1](#).

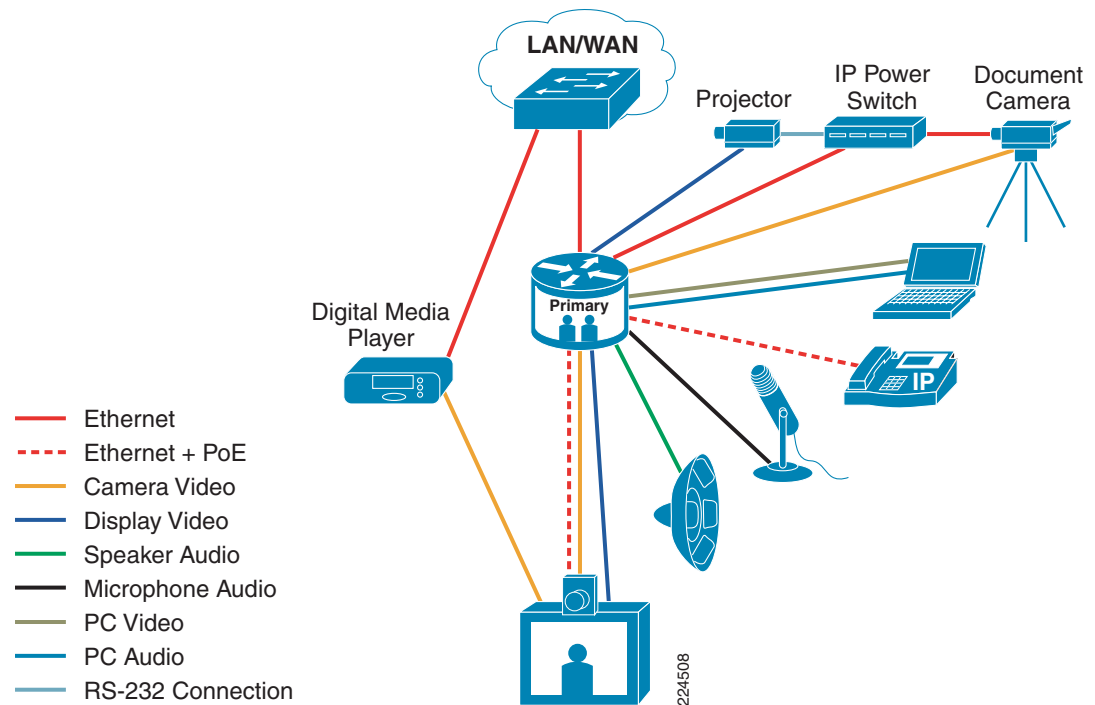
Figure 2-1 Connectivity Schematic for a CTS-1000 System



Connecting a CTS-500 System

The CTS-500 is an enclosed system that includes a 37" display with integrated codec, camera, microphone, and speaker. Its only connection point is an RJ-45 cable which connects it to the network. For purposes of this discussion, the CTS-500 is similar to the CTS-1000 in terms of connectivity, but has a connection for an optional digital media player. This can be used to display content when the CTS-500 is not being used for TelePresence Meetings. The connectivity of the CTS-500 is shown in Figure 2-2.

Figure 2-2 Connectivity Schematic for a CTS-500 System



Connecting a CTS-3000 System

The CTS-3000 system includes:

- One Cisco TelePresence primary codec
- Two Cisco TelePresence secondary codecs
- One Cisco Unified 7975G IP phone
- Three 65" plasma displays
- Three high-definition cameras
- Three microphones
- Three speakers
- One input for auxiliary audio
- One input for auxiliary video which can be used for a document camera or PC

As with the CTS-1000 system, the primary codec is the central part of the CTS-3000 system to which all other components interconnect.

Specifically, the Cisco Unified 7975G IP phone connects to the TelePresence primary codec via an RJ-45 cable that provides it network connectivity and 802.3af Power-over-Ethernet (PoE).

A video cable connects the primary codec to the center 65” plasma display; another of these cables connects the right display to the (right) secondary codec, and a third connects the left display to the (left) secondary codec. As with the CTS-1000 system, this cable is essentially an HDMI cable, but with a proprietary element for carrying management information instead of audio signals (as the audio signals are processed independently by the master codec). Each of these secondary codecs, in turn, are connected to the primary codec via a RJ-45 cable; however, no 802.3af PoE is required over these Ethernet links as the secondary codecs have independent power supplies.

Three cameras are mounted on the central display and each camera is connected to its respective codec:

- The left camera is connected to the (left) secondary codec.
- The center camera is connected to the primary codec.
- The right camera is connected to the (right) secondary codec.

Each camera connects to its respective codec via two cables: a RJ-45 cable, which provides 802.3af PoE and network connectivity to the camera and a video cable to carry the video signals to the codec.

Additionally, three speaker cables and three microphone cables connect the (left, center, and right) speakers and (left, center, and right) microphones to the primary codec, respectively.

The primary codec also has inputs for auxiliary audio and auxiliary video. Auxiliary video can come from a PC connection or from a document camera connection. An IP power switch (IPS) provides control for the on/off function of the document camera, attached projector, as well as the lighting shroud of the CTS unit via an Ethernet connection.

Finally, an RJ-45 cable provides 10/100/1000 Ethernet connectivity from the primary codec to the network infrastructure. These interconnections for a CTS-3000 system are illustrated in [Figure 2-3](#).

Specifically, the Cisco Unified 7975G IP phone connects to the TelePresence primary codec via an RJ-45 cable that provides it with network connectivity and 802.3af Power-over-Ethernet (PoE).

A video cable connects the primary codec to the center 65" plasma display, another cable connects the right display to the (right) secondary codec, and a third connects the left display to the (left) secondary codec. As with the CTS-3000 system, this cable is essentially an HDMI cable, but with a proprietary element for carrying management information instead of audio signals (as the audio signals are processed independently by the master codec). Each of these secondary codecs, in turn, are connected to the primary codec via an RJ-45 cable; however, no 802.3af PoE is required over these Ethernet links as the secondary codecs have independent power supplies.

Three cameras are mounted on the central display and each camera is connected to its respective codec:

- The left camera is connected to the (left) secondary codec.
- The center camera is connected to the primary codec.
- The right camera is connected to the (right) secondary codec.

Each camera connects to its respective codec via two cables:

- RJ-45 cable, which provides 802.3af PoE and network connectivity to the camera
- Video cable to carry the video signals to the codec

Additionally, three speaker cables connect the (left, center, and right) speakers to the primary codec, respectively.

One microphone cable connects the center microphone to the primary codec. The remaining eight microphones are connected to the audio extension box, which is in turn connected to the primary codec. The audio extension box also houses the HDMI splitter. The HDMI splitter connects to the auxiliary video output of the primary codec. Up to four displays or a projector and three displays can be connected to the HDMI ports on the audio extension box.

The primary codec also has inputs for low-speed (5 frames per second) auxiliary audio and auxiliary video inputs. Video input can come from a PC or optional document camera. An IP power switch (IPS) provides control for the on/off function of the document camera, attached projector, as well as the lighting shroud of the CTS unit via an Ethernet connection. Optionally, another secondary codec can be connected to the primary codec to provide high-speed (30 frames per second) auxiliary video input. The auxiliary codec is connected to the primary codec via the RJ-45 cable from the Ethernet port normally used for the document camera. Auxiliary audio is still connected to the primary codec.

A connection for an optional headset is also provided on the primary codec of the CTS-3200.

Finally, an RJ-45 cable provides 10/100/1000 Ethernet connectivity from the primary codec to the network infrastructure. These interconnections for a CTS-3200 system are illustrated in [Figure 2-4](#) and [Figure 2-5](#).

Figure 2-4 Connectivity Schematic for a CTS-3200 System with Low-Speed Auxiliary Input

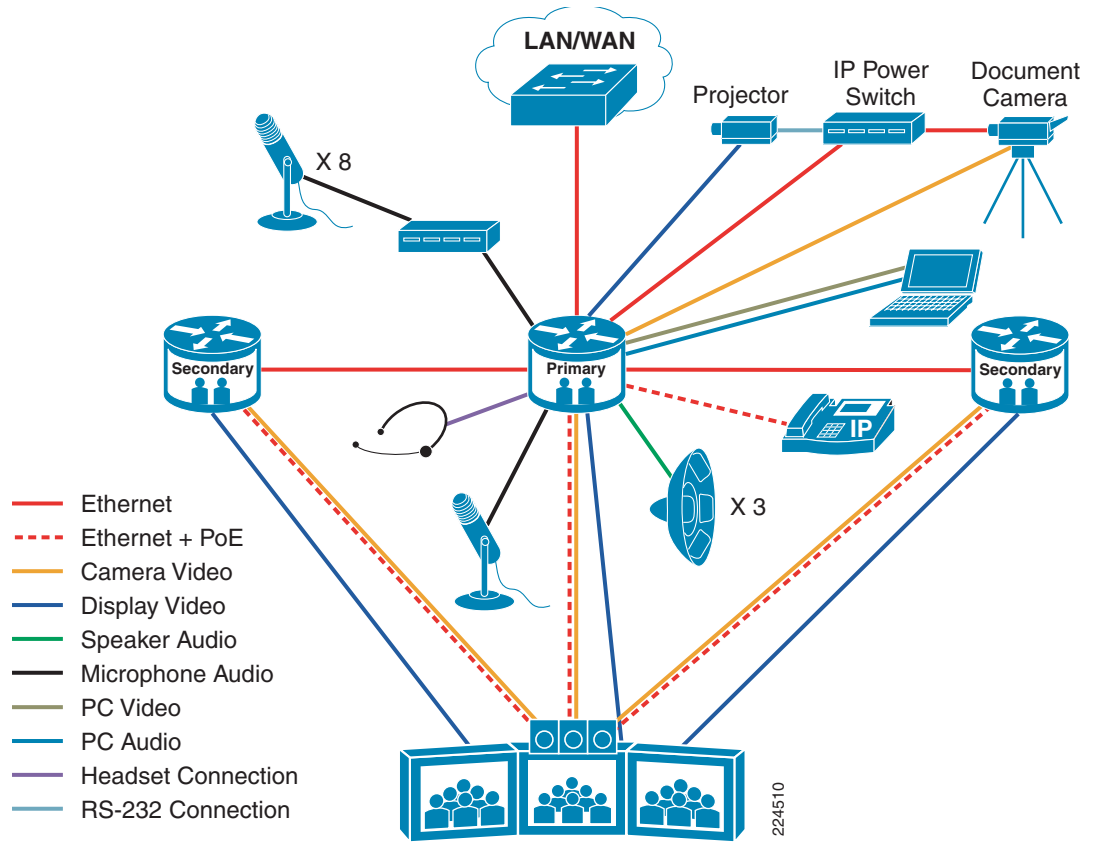
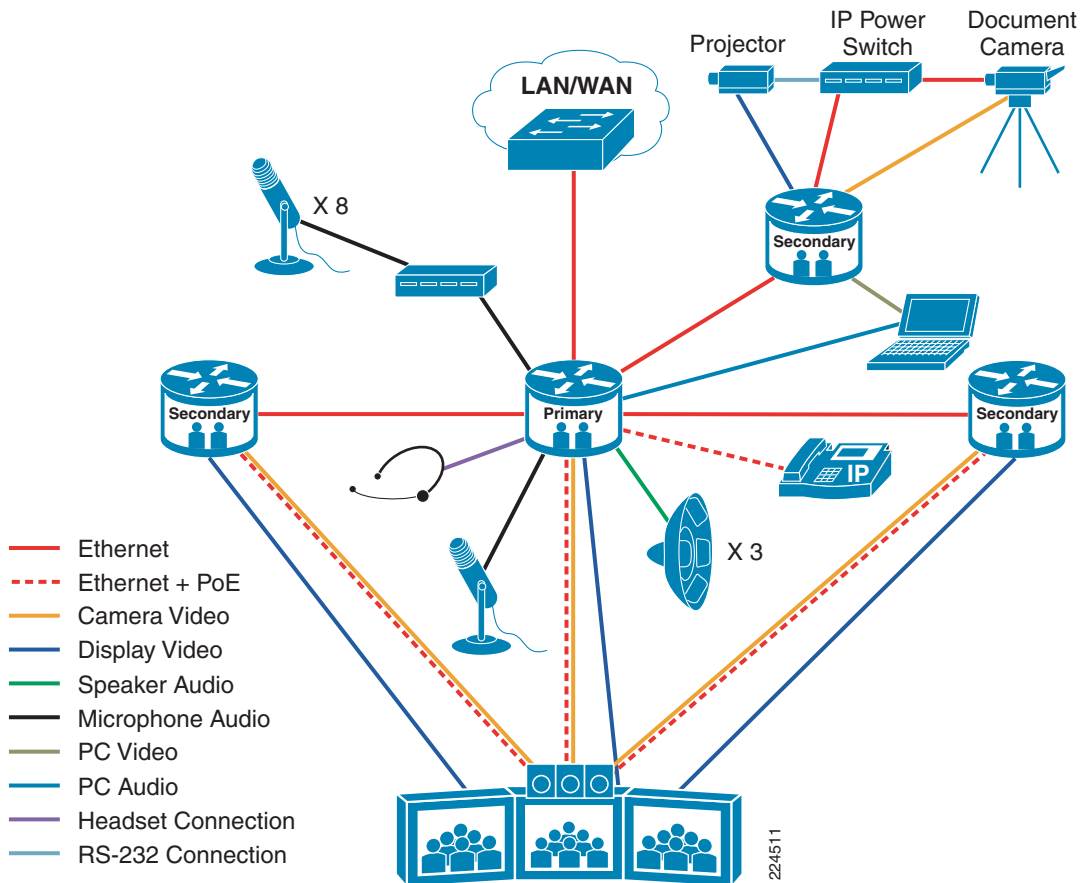


Figure 2-5 Connectivity Schematic for a CTS-3200 System with High-Speed Auxiliary Input

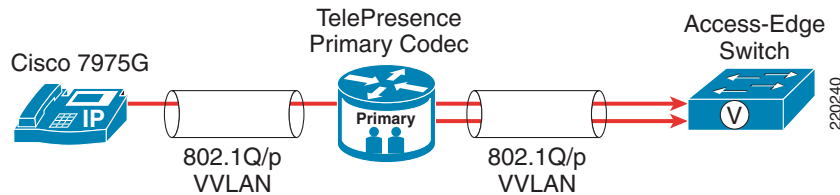


Cisco TelePresence Network Interaction

The primary codec is the interface between the CTS endpoint system and the network infrastructure. The primary codec connects to the network access edge switch via a RJ-45 10/100/1000 port. The access edge Catalyst switch that it connects to provides IP services, 802.1Q/p VLAN services, QoS services, and security services to the TelePresence endpoint.

Additionally, the primary codec provides a RJ-45 connection to the Cisco Unified 7975G IP phone, to which it supplies 802.3af PoE. When the IP phone boots up, it sends a Cisco Discovery Protocol (CDP) message to the primary codec. The codec receives this CDP message and passes it on to the access edge switch, supplementing it with its own CDP advertisement. The access edge switch and codec exchange CDP messages and the switch (if configured according to best practice recommendations for IP telephony deployments) places the primary codec and the 7975G IP phone in a 802.1Q Voice VLAN (VVLAN), wherein 802.1Q/p Class of Service (CoS) markings are trusted. The primary Codec passes 802.1Q tags between the 7975G IP phone and the network access edge switch, extending the VVLAN all the way to the IP phone. This 802.1Q/p VVLAN assignment is illustrated in [Figure 2-6](#).

Figure 2-6 Voice VLAN Extension Through Cisco TelePresence Primary Codec



Note

The above network interaction assumes that CDP is enabled and Voice VLANs are configured. If this is not the case, then the network interaction begins with the DHCP requests described next.

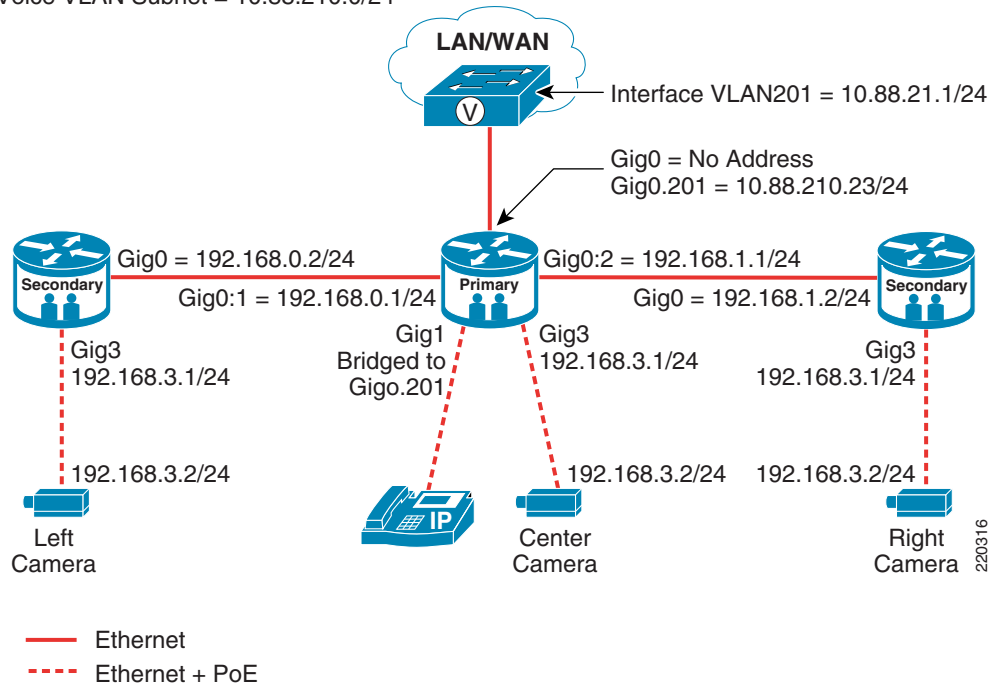
If configured for dynamic IP addressing, the 7975G IP phone and the primary codec each generate a Dynamic Host Configuration Protocol (DHCP) request to the network and are supplied with IP addresses (one for the IP phone and another for the primary codec). The DHCP server may also provide the IP phone and primary codec the IP address of the download server, via DHCP option 150, from which they download their configuration files and firmware loads. This function is often provided by the Cisco Unified Communications Manager (CUCM) server. Alternatively, either or both of the devices may be configured with a static IP address and TFTP server address.

Additionally, it is important to note that the TelePresence systems utilize a private network for internal communications between the primary and secondary codecs, as well as between codecs and cameras. By default the internal address range used is 192.168.0.0/24 through 192.168.4.0/24; however, if the TelePresence codec receives a 192.168.x.x address from the network, then the internal private network will switch to 10.0.0.0/24 through 10.0.4.0/24. A default internal network IP address assignment is illustrated in [Figure 2-7](#).

Figure 2-7 Default TelePresence Internal IP Addressing Scheme**Example:**

Voice VLAN ID = 201

Voice VLAN Subnet = 10.88.210.0/24

**Note**

Even though only 192.168.0.0/24 through 192.168.3.0/24 are illustrated in [Figure 2-7](#), 192.168.4.0/24 is reserved within the system for future (internal) use.

Similarly, if the TelePresence system is using 10.0.0.0/24 through 10.0.3.0/24 for its internal networking address range, then 10.0.4.0/24 is reserved within the system for future (internal) use.

It is important to note three key points regarding the internal networking of TelePresence systems:

- From the network's perspective, the TelePresence primary codec appears as a single endpoint device with a single IP address (but remember, the 7975G IP Phone also appears as a separate endpoint device with its own IP address).
- The internal components (such as secondary codecs and cameras) do not receive a default gateway. Therefore, they cannot route beyond the primary codec.
- If the primary codec is using 192.168.0.0/24 through 192.168.4.0/24 as its internal networking addresses (which is the default), then it is not able to connect to external servers or endpoints that are using these same addresses (as it will attempt to reach such addresses via its internal network, not its external default gateway). Conversely, if the primary codec has been assigned an IP address from the network in the 192.168.x.x range, then it uses internal networking addresses in the range of 10.0.0.0/24 through 10.0.4.0/24, and similarly, is not able to connect to external servers or endpoints that may be using these same addresses. [Table 2-1](#) summarizes the IP addressing best practices for networks supporting TelePresence.

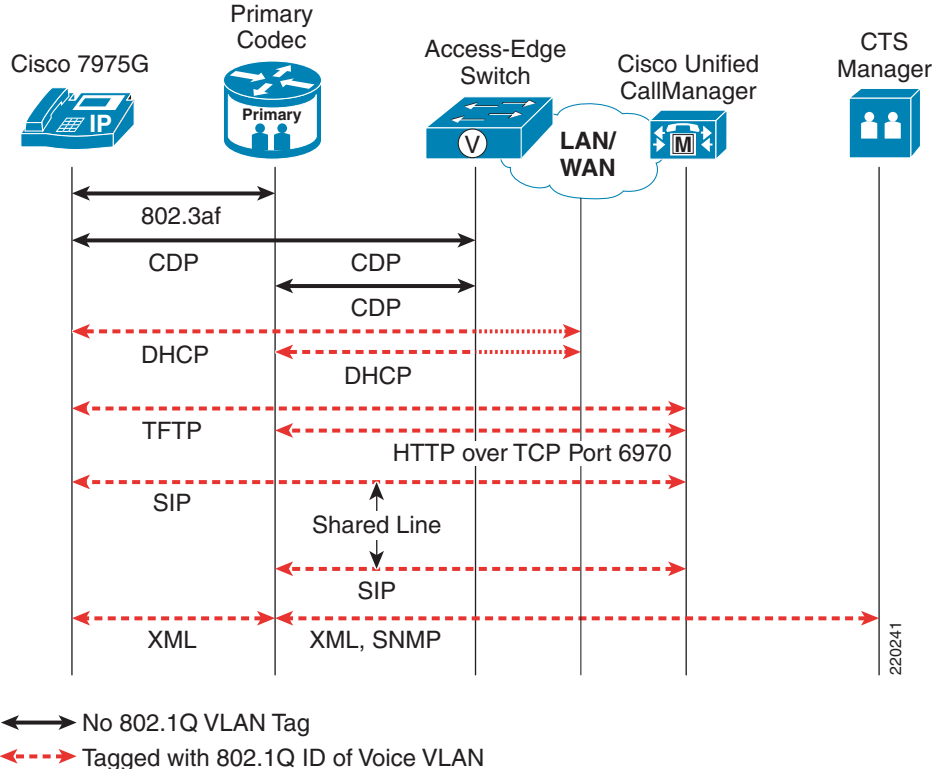
Table 2-1 *TelePresence Network IP Addressing Best Practices*

For Environments Where the CTS Uses 192.168.x.x for its Internal Communications. Avoid Using the Following Subnets:	For Environments Where the CTS Uses 10.x.x.x for its Internal Communications. Avoid Using the Following Subnets:
192.168.0.0/24	10.0.0.0/24
192.168.1.0.24	10.0.1.0/24
192.168.2.0.24	10.0.2.0/24
192.168.3.0.24	10.0.3.0/24
192.168.4.0.24	10.0.4.0/24

Provided there are no IP addressing issues, as described above, the IP phone then initiates a Trivial File Transfer Protocol (TFTP) session with the Cisco Unified Communications Manager (CUCM) to download its configuration and firmware files. The primary codec initiates an HTTP session over TCP port 6970 for its configuration and firmware files. Note that DNS may also be required to translate the CUCM hostname to an IP address.

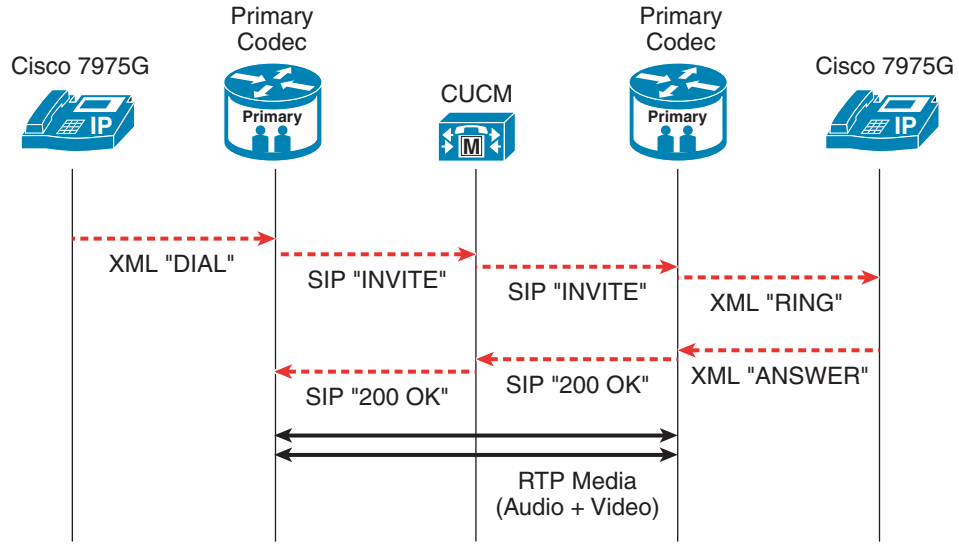
The primary codec then communicates with CUCM via Session Initiation Protocol (SIP). The Cisco 7975G IP Phone also communicates with CUCM via SIP, identifying itself as a shared line with the primary codec. Additional messaging occurs between the 7975G IP phone, the TelePresence primary codec, and the Cisco TelePresence Manager via Extensible Markup Language (XML), as well as Simple Network Management Protocol (SNMP). These network protocol interactions are illustrated in [Figure 2-8](#).

Figure 2-8 Cisco TelePresence Network Control, Management, and Signaling Protocols



Once the TelePresence system has completed these protocol interactions, it is ready to place and receive calls. When a call is initiated, the Cisco 7975G IP phone sends an XML Dial message to its primary codec, which forwards the request as a SIP Invite message to the Cisco Unified Communications Manager. CUCM, in turn, forwards the SIP Invite message to the destination TelePresence codec, which forwards the message as an XML Ring message to its 7975G associated IP phone. The TelePresence primary codec can be set to automatically answer the incoming call or can be set to send an incoming call alert to the 7975G IP phone. If set to auto-answer, the codec answers the call immediately and sends a SIP OK message to CUCM. If auto-answer is not enabled, when the user presses the Answer softkey on the 7975G IP phone, the 7975G IP phone replies with a XML Answer message to the receiving TelePresence primary codec, and the codec in turn sends a SIP 200 OK message to CUCM. CUCM relays this SIP 200 OK message to the originating TelePresence primary codec and the call is established. Real-time media, both audio and video, is then passed between the TelePresence primary codecs over Real Time Protocol (RTP). The signaling and media paths for Cisco TelePresence are illustrated in [Figure 2-9](#).

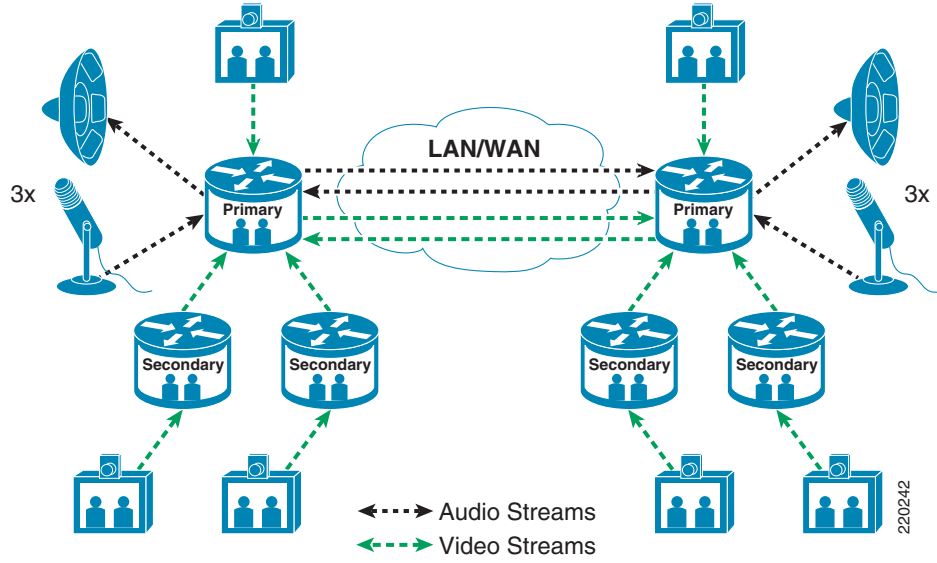
Figure 2-9 Cisco TelePresence Signaling and Media Paths



- - - - - 7975G Signaling Note: Signaling has been simplified for the purpose of this figure. 220213
< - - - - - > Media

CTS-1000 and CTS-500 systems send only one audio and one video stream (excluding auxiliary audio and video inputs for the moment). On the other hand, CTS-3000 and CTS-3200 primary codecs process three separate audio and three separate video streams. However, these codecs do not send three separate audio streams and three separate video streams over the network. Rather, CTS-3000 and CTS-3200 primary codecs multiplex the three audio streams into one RTP stream and three video streams into one RTP stream, and hence send only a single audio and a single video stream over the network. These streams, in turn, are de-multiplexed by the receiving codec. The multiplexing of audio and video streams performed by the CTS-3000 primary codecs is illustrated in [Figure 2-10](#). Auxiliary audio and video inputs are also multiplexed into the same audio and video streams. Therefore, in the case of the CTS-1000 or CTS-500, the primary video and auxiliary video are multiplexed into one outgoing video stream; likewise the primary audio and auxiliary audio are multiplexed into one outgoing audio stream. In the case of the CTS-3000 or CTS-3200, the auxiliary video is treated as the 4th video channel and multiplexed in with the rest of the video; likewise the auxiliary audio is treated as the 4th audio channel and multiplexed in with the rest of the audio.

Figure 2-10 CTS-3000 Multiplexing of Audio and Video Streams





CHAPTER 3

TelePresence Network Deployment Models

Introduction

TelePresence systems can be deployed over enterprise networks in one of four principle ways:

- [Intra-Campus Deployment Model](#)
- [Intra-Enterprise Deployment Model](#)
- [MultiPoint Deployment Model](#) (see [Point-to-Point versus Multipoint](#))
- [Inter-Enterprise/Business-to-Business Deployment Model](#)

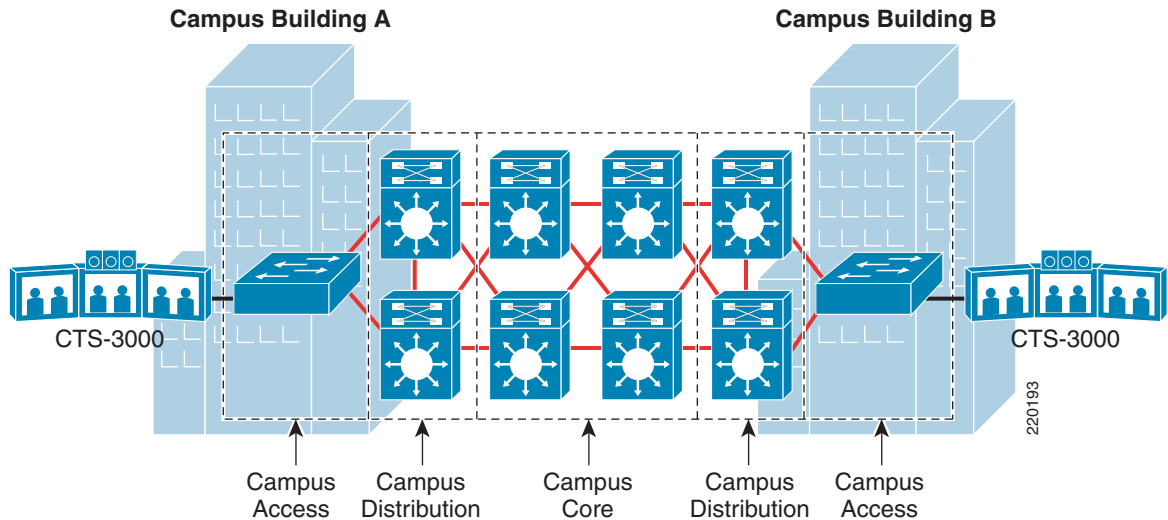
The following sections provide an overview of these TelePresence network deployment models, as well as logical phases of TelePresence deployments. In comparison, CUCM deployment models are discussed in detail in [Chapter 9, “Call Processing Deployment Models.”](#)

Intra-Campus Deployment Model

The intra-campus network deployment model has TelePresence systems limited to a single enterprise campus or between sites interconnected via a high-speed (1 Gigabit or higher) Metropolitan Area Network (MAN). This deployment model is applicable for enterprises that have a large number of buildings within a given campus and employees who are often required to drive to several different buildings during the course of the day to attend meetings. Deploying multiple TelePresence systems intra-campus can reduce time lost by employees driving between buildings to attend meetings, without sacrificing meeting effectiveness, and thus improve overall productivity. The intra-campus deployment model is also commonly used in conjunction with the other two: where customers deploy multiple CTS rooms within their headquarters campus to meet demand for room availability as part of a global intra-enterprise or inter-enterprise deployment.

The network infrastructure of an intra-campus deployment model is predominantly Cisco Catalyst switches connecting via GigE or 10GigE links. The intra-campus TelePresence deployment model is illustrated in [Figure 3-1](#).

Figure 3-1 *TelePresence Intra-Campus Network Deployment Model*



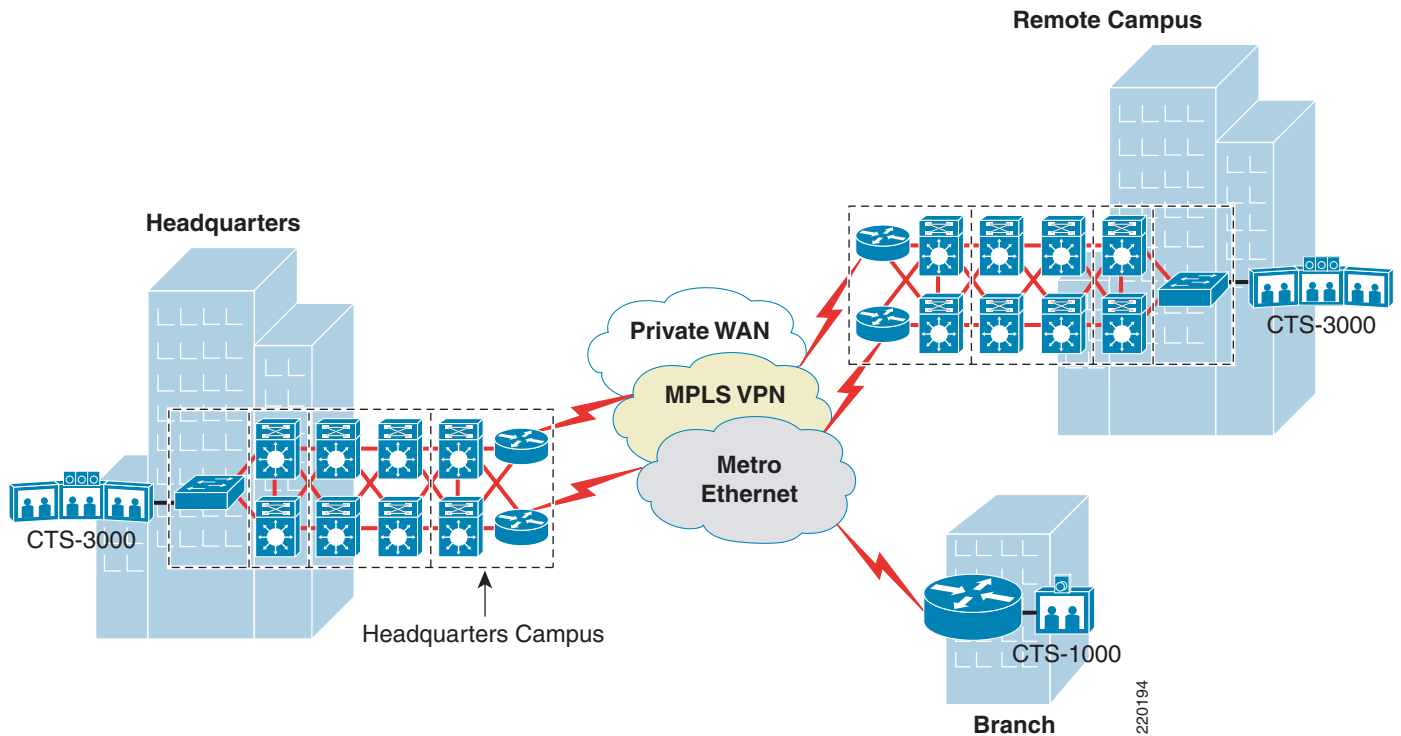
Intra-Enterprise Deployment Model

The intra-enterprise network deployment model for TelePresence systems connects not only buildings within a campus, but also geographically-separated campus sites and branch offices. The intra-enterprise model expands on the intra-campus model to include sites connected via a Wide Area Network (< 1 Gigabit).

The intra-enterprise deployment model is suitable for businesses that often require employees to travel extensively for internal meetings. Deploying TelePresence systems within the enterprise not only improves productivity—by saving travel time—but also reduces travel expenses. Furthermore, the overall quality of work/life is often improved when employees have to travel less.

The network infrastructure of an intra-enterprise deployment model is a combination of Cisco Catalyst switches within the campus and Cisco routers over the WAN, which may include private WANs, MPLS VPNs, or Metro Ethernet networks. WAN speeds may range from 34 Mbps E3 circuits to 1 Gbps OC-192 circuits. The intra-enterprise TelePresence deployment model is illustrated in [Figure 3-2](#).

Figure 3-2 TelePresence Intra-Enterprise Network Deployment Model



Cisco Powered Networks

A valuable consideration when selecting WAN/VPN service providers is to identify those that have achieved Cisco Powered Network designation. These providers have earned the Cisco Powered designation by maintaining high levels of network quality and by basing their WAN/VPN services end-to-end on Cisco equipment.

In addition, an increasing number of Cisco Powered providers have earned the QoS Certification for WAN/VPN services. This means that they have been assessed by a third party for the ability of their SLAs to support real-time voice and video traffic, and for their use of Cisco best practices for QoS. For a list of recommended service providers, see the following URL: <http://www.cisco.com/cpn>.

The use of Cisco Powered networks is recommended—but not mandatory—for Cisco TelePresence intra-enterprise deployments. The key is meeting the service levels required by TelePresence, which are detailed in Chapter 4, “Quality of Service Design for TelePresence.”

Point-to-Point versus Multipoint

In both the intra-campus and inter-enterprise deployment models, customers may also deploy multipoint TelePresence resources to facilitate multi-site meetings (meetings with three or more TelePresence rooms). These resources may be located at any one of the campus locations or may be located within the service provider cloud as either a co-located resource or a managed/hosted resource.

Multipoint platforms and network design recommendations, such as additional bandwidth and latency considerations, Cisco TelePresence Multipoint switch considerations, scaling considerations, etc., are discussed in further detail in [Chapter 10, “Cisco TelePresence Multipoint Solution Essentials,”](#) [Chapter 11, “Cisco Multipoint Technology and Design Details,”](#) and [Chapter 12, “Cisco TelePresence Multipoint Solution Circuit and Platform Recommendations.”](#)

Inter-Enterprise/Business-to-Business Deployment Model

The inter-enterprise network deployment model connects not only TelePresence systems within an enterprise, but also allows for TelePresence systems within one enterprise to call systems within another enterprise. The inter-enterprise model expands on the intra-campus and intra-enterprise models to include connectivity between different enterprises. This is also referred to as the business-to-business (B2B) TelePresence deployment model.

The inter-enterprise model offers the most flexibility and is suitable for businesses that often require employees to travel extensively for both internal and external meetings. In addition to the business advantages of the intra-enterprise model, the B2B TelePresence deployment model lets employees maintain high-quality customer relations, without the associated costs of travel time and expense.

The network infrastructure of the inter-enterprise/B2B deployment model builds on the intra-enterprise model and requires the enterprises to share a common MPLS VPN service provider (SP). Additionally, the MPLS VPN SP must have a “shared services” Virtual Routing and Forwarding (VRF) instance provisioned with a Cisco IOS XR Session/Border Controller (SBC).

The Cisco SBC bridges a connection between two separate MPLS VPNs to perform secure inter-VPN communication between enterprises. Additionally, the SBC provides topology and address hiding services, NAT and firewall traversal, fraud and theft of service prevention, DDoS detection and prevention, call admission control policy enforcement and guaranteed QoS.

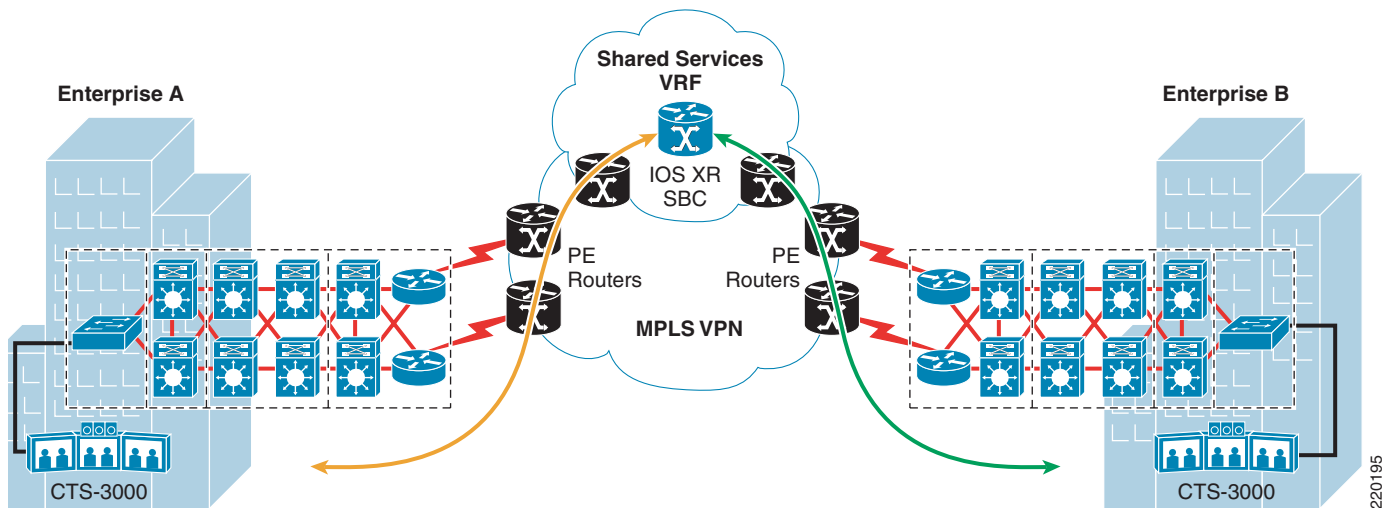


Note

For more information about Cisco IOS XR SBC functionality and deployment models, refer to: http://www.cisco.com/univercd/cc/td/doc/product/ioxsoft/iox34/cgcr34/sbc_c34/sbc34abt.htm

The inter-enterprise/B2B TelePresence deployment model is illustrated in [Figure 3-3](#).

Figure 3-3 TelePresence Inter-Enterprise Network Deployment Model



The initial release of the B2B solution requires a single SP to provide the shared services to enterprise customers, which includes the secure bridging of customer MPLS VPNs. However, as this solution evolves, multiple providers will be able to peer and provide B2B services between them, which will no longer require that both enterprise customers share the same SP.

Hosting and Management Options

While the focus of this paper is TelePresence deployments within the enterprise, several of these options could be hosted or managed by SPs. For example, the Cisco Unified Communications Manager (CUCM) and Cisco TelePresence Manager (CTS-MAN) servers and multipoint resources may be located on-premise at one of the customer campus locations, co-located within the SP network (managed by the enterprise) or hosted within the SP network (managed by the SP). However, with the exception of inter-VPN elements required by providers offering B2B TelePresence services, the TelePresence solution components and network designs remain fundamentally the same whether the TelePresence systems are hosted/managed by the enterprise or the SP.

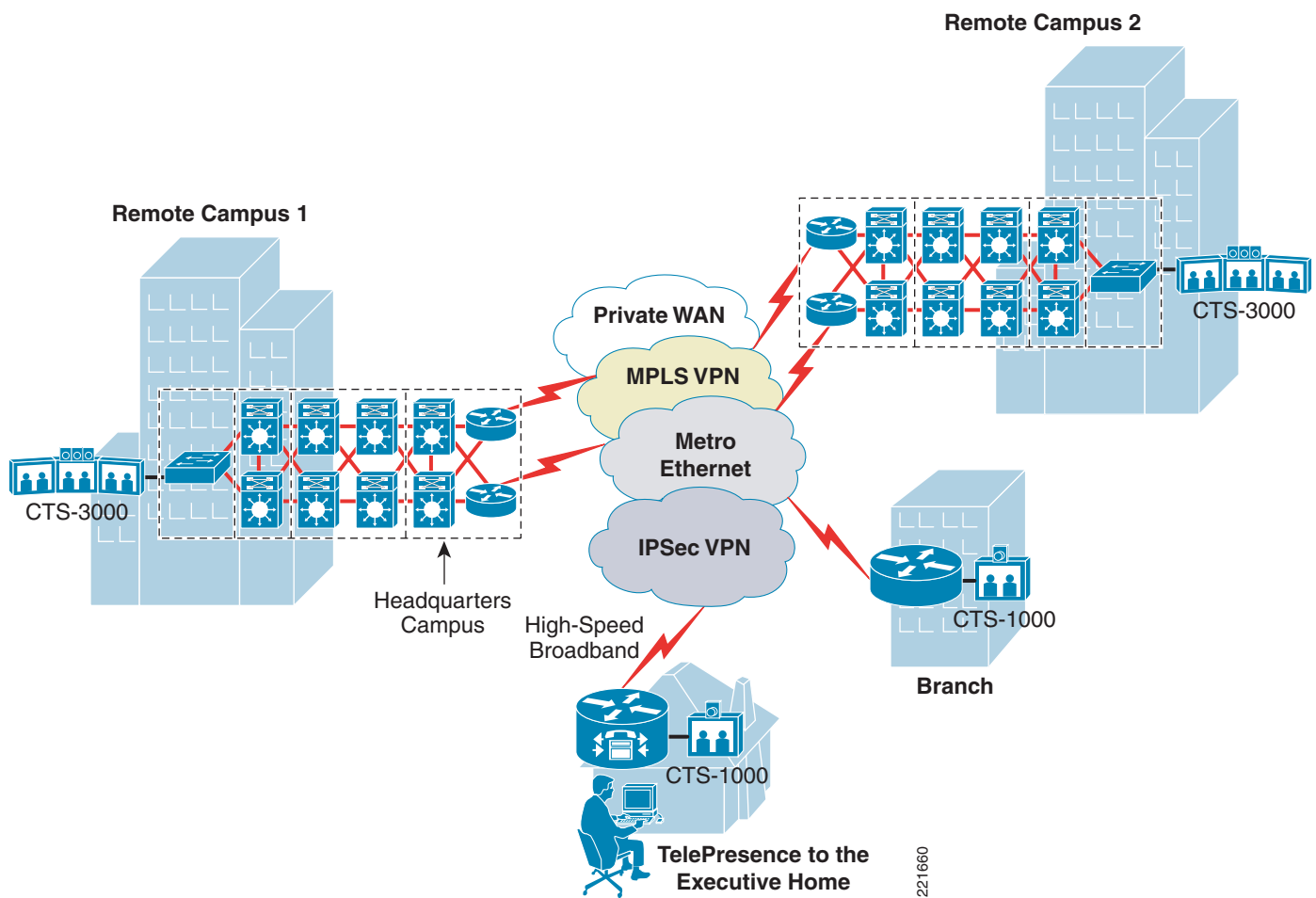
TelePresence Phases of Deployment

As TelePresence technologies evolve, so too will the complexity of deployment solutions. Therefore, enterprise customers will likely approach their TelePresence deployments in phases, with the main phases of deployment being:

- Phase 1. Intra-Campus/Intra-Enterprise Deployments—Most enterprise customers will likely begin their TelePresence rollouts by provisioning (Point-to-Point) Intra-Enterprise TelePresence deployments. This model could be viewed as the basic TelePresence building block, on which more complex models may be added.
- Phase 2. Intra-Enterprise MultiPoint Deployments—As collaboration requirements may not always be facilitated with Point-to-Point models, the next logical phase of TelePresence deployment would be to introduce multipoint resources to the Intra-Enterprise deployment model. Phases 1 and 2 may often be undertaken simultaneously.

- Phase 3. Business-to-Business Deployments—To expand the application and business benefits of TelePresence meetings to include external (customer- or partner-facing) meetings, a Business-to-Business deployment model can be subsequently overlaid on top of either a Point-to-Point or a MultiPoint Intra-Enterprise deployment.
- Phase 4. TelePresence to the Executive Home—Due to the high executive-perk appeal of TelePresence and the availability of high-speed residential bandwidth options (such as fiber to the home), some executives may benefit greatly from deploying TelePresence units to their residences. Technically, this is simply an extension of the Intra-Enterprise model, but for the purposes of this document it is viewed as a separate phase due to the unique provisioning and security requirements posed by such residential TelePresence deployments.

Figure 3-4 *TelePresence to the Executive Home (an Extension of the Intra-Enterprise Deployment Model)*





CHAPTER 4

Quality of Service Design for TelePresence

Overview

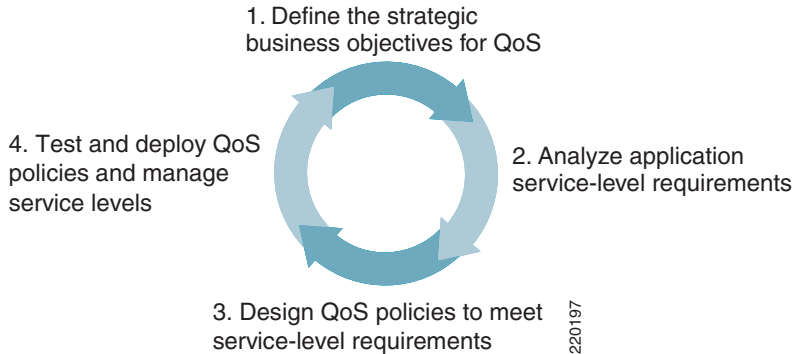
A major benefit of Cisco's TelePresence solution over competitive offerings is that the realtime, high-definition video and audio are transported over a converged IP network rather than a dedicated network (although dedicated networks are also supported). The key enabling technology to accomplish this convergence is Quality of Service (QoS).

QoS technologies refer to the set of tools and techniques to manage network resources, such as bandwidth, latency, jitter, and loss. QoS technologies allow different types of traffic to intelligently contend for network resources. For example, voice and realtime video—such as TelePresence—may be granted strict priority service, while some critical data applications may receive (non-priority) preferential services and some undesired applications may be assigned deferential levels of service. Therefore, QoS is a critical, intrinsic element for the successful network convergence of voice, video, and data.

There are four principal phases to a successful QoS deployment:

- Clearly define the strategic business objectives of the QoS deployment.
- Analyze application service-level requirements.
- Design (and test) QoS policies to accommodate service level requirements.
- Roll out the QoS policies and monitor service levels.

These phases are sequential and the success of each subsequent phase directly depends on how well the previous phase has been addressed. Furthermore, the entire process is generally cyclical, as business applications and objectives evolve over time and their related QoS policies periodically need to be adjusted to accommodate (see [Figure 4-1](#)).

Figure 4-1 The Four Phases of Successful QoS Deployments

The following sections examine how each of these phases relate to a successful deployment of QoS for TelePresence.

Defining the Strategic Business Objective for QoS for TelePresence

QoS technologies are the enablers for business/organizational objectives. Therefore, the way to begin a QoS deployment is not to activate QoS features simply because they exist, but to start by clearly defining the QoS-related business objectives of the organization.

For example, among the first questions that arise during a QoS deployment are: How many traffic classes should be provisioned for? And what should they be? To help answer these fundamental questions, QoS-related organizational objectives need to be defined, such as:

- Is the business objective to enable TelePresence only? Or is VoIP also required to run over the converged network?
- Are there any non-realtime applications that are considered critical to the core business objectives? If so, what are they?
- Are there applications which should be squelched (i.e., deferential treatment)? If so, what are they?

The answers to these questions define the applications that require QoS policies, either preferential QoS or deferential QoS. Each application that has a unique service level requirement—whether preferential or deferential—requires a dedicated service class to deliver and guarantee the requisite service levels.

Additionally, Cisco offers a non-technical recommendation for this first phase of a successful QoS deployment, namely to always seek executive endorsement of the QoS business objectives prior to design and deployment. This is because QoS is a system of managed application preference and as such often includes political and organizational repercussions when implemented. To minimize the effects of these non-technical obstacles to deployment, it is recommended to address these political and organizational issues as early as possible, garnishing executive endorsement whenever possible.

Analyzing the Service Level Requirements of TelePresence

Once the applications requiring QoS have been defined by the organization business objectives, then the network administrators must carefully analyze the specifics of the service levels required by each application to be able to define the QoS policies to meet them. The service level requirements of realtime applications, such as TelePresence, are defined by the following four parameters:

- Bandwidth
- Latency (delay)
- Jitter (variations in delay)
- Packet loss

TelePresence Bandwidth Requirements

Cisco TelePresence systems are currently available in one screen (CTS-1000) and three screen (CTS-3000) configurations. A CTS-3000 obviously has greater bandwidth requirements than a CTS-1000, but not necessarily by a full-factor of three, as will be shown. Furthermore, the resolution of each CTS-1000 or CTS-3000 system can be set to 720p or 1080p (full HDTV); the resolution setting also significantly impacts the bandwidth requirements of the deployed TelePresence solution.

As discussed in [Chapter 1, “Cisco TelePresence Solution Overview,”](#) Cisco TelePresence has even more levels of granularity in overall image quality within a given resolution setting, as the motion handling quality can also be selected. Therefore, TelePresence supports three levels of motion handling quality within a given resolution, specifically 720p-Good, 720p-Better, and 720p-Best, as well as 1080p-Good, 1080p-Better, and 1080p-Best. Each of these levels of resolution and motion handling quality results in slightly different bandwidth requirements, as detailed in [Table 4-1](#).

To keep the following sections and examples simple to understand, only two cases will be broken down for detailed analysis: 720p-Good and 1080p-Best.

Let’s break down the bandwidth requirements of the maximum bandwidth required by a CTS-1000 system running at 720p-Good, with an auxiliary video stream (a 5 frame-per-second video channel for sharing Microsoft PowerPoint or other collateral via the data-projector) and an auxiliary audio stream (for at least one additional person conferenced in by an audio-only bridge). The bandwidth requirements by component are:

1 primary video streams @ 1 Mbps:	1,000 Mbps (1 Mbps)
1 primary audio streams @ 64 Kbps:	64 Kbps
1 auxiliary video stream (5 fps):	500 Kbps
1 auxiliary audio stream:	<u>64 Kbps</u>
Total audio and video bandwidth (not including burst and network overhead):	1,628 Kbps (1.628 Mbps)

The total bandwidth requirements—without network overhead—of such a scenario would be 1.628 Mbps. However a 10% burst factor on the video channel, along with the IP/UDP/RTP overhead (which combined amounts to 40 bytes per packet) must also be taken into account and provisioned for, as must media-specific Layer 2 overhead. In general, video—unlike voice—does not have clean formulas for calculating network overhead because video packet sizes and rates vary proportionally to the degree of

motion within the video image itself. From a network administrator's point of view, bandwidth is always provisioned at Layer 2, but the variability in the packet sizes and the variety of Layer 2 mediums the packets may traverse from end-to-end make it difficult to calculate the real bandwidth that should be provisioned at Layer 2. Cisco TelePresence video packets average 1,100 bytes per packet. However, the conservative rule of thumb that has been thoroughly tested and widely deployed is to overprovision video bandwidth by 20%. This accommodates the 10% burst and the Layer 2-Layer 4 network overhead.

With this 20% overprovisioning rule applied, the requisite bandwidth for a CTS-1000 running at 720p-Good becomes 2 Mbps (rounded).

Now, let's break down the maximum bandwidth required by a CTS-3000 system running at full 1080p-Best, with an auxiliary video stream and an auxiliary audio stream. The detailed bandwidth requirements are:

3 primary video streams @ 4 Mbps each:	12,000 Kbps (12 Mbps)
3 primary audio streams @ 64 Kbps each:	192 Kbps
1 auxiliary video stream:	500 Kbps
1 auxiliary audio stream:	<u>64 Kbps</u>
Total audio and video bandwidth (not including burst and network overhead):	12,756 Kbps (12.756 Mbps)

With the 20% overprovisioning rule applied, the requisite bandwidth for a CTS-3000 running at 1080p-Best becomes approximately 15 Mbps (with a bit of rounding applied). This value of 15 Mbps for a CTS-3000 at 1080p-Best is used in most of the examples throughout this design guide.

It is important to note that as the Cisco TelePresence software continues to evolve and add new feature support, the bandwidth requirements for TelePresence will correspondingly evolve and expand. For example, [Table 4-1](#) shows the bandwidth requirements for CTS software version 1.2, with and without network overhead, of CTS-1000 and CTS-3000 systems running at 720p and 1080p with all grades of motion handling quality (Good, Better, and Best).

Table 4-1 Cisco TelePresence Software Version 1.2 Bandwidth Requirements

Resolution	1080p	1080p	1080p	720p	720p	720p
Motion Handling	Best	Better	Good	Best	Better	Good
Video per Screen (kbps)	4000	3500	3000	2250	1500	1000
Audio per Microphone (kbps)	64	64	64	64	64	64
(5 fps) Auto Collaborate Video Channel (kbps)	500	500	500	500	500	500
Auto Collaborate Audio Channel (kbps)	64	64	64	64	64	64
CTS-500/1000 Total Audio and Video (kbps)	4,628 ¹	4,128 ¹	3,628 ¹	2,878 ¹	2,128 ¹	1,628 ¹
CTS-3000/3200 Total Audio and Video (kbps)	12,756	11,256	9,756	7,506	5,256	3,756
CTS-500/1000 total bandwidth (Including Layer 2-Layer 4 overhead)	5.5 Mbps ¹	5.0 Mbps ¹	4.3 Mbps ¹	3.4 Mbps ¹	2.5 Mbps ¹	2 Mbps ¹
CTS-3000/3200 total bandwidth (Including Layer 2-Layer 4 overhead)	15.3 Mbps	13.5 Mbps	11.7 Mbps	9.0 Mbps	6.3 Mbps	4.5 Mbps

1. The CTS-1000 transmits up to 128kbps of audio, but can receive up to 256kbps when participating in a meeting with a CTS-3000.

**Note**

These bandwidth numbers represent the worst-case scenarios (i.e., peak bandwidth transmitted during periods of maximum motion within the encoded video). Normal use (i.e., average bandwidth), with users sitting and talking and gesturing naturally, typically generates only about 60-80% of these maximum bandwidth rates. This means that a CTS-3000 running at 1080-Best averages only 10-12 Mbps and a CTS-1000 running at 720-Good averages only 1.2-1.6 Mbps.

Release 1.3 of CTS software introduced support for an interoperability feature, which allows for TelePresence systems to interoperate with H.323-based video-conferencing systems. From a bandwidth perspective, the only change to existing TelePresence flows is that there is an additional video channel transmitted (768 Kbps) and an additional audio channel transmitted (64 Kbps) by the TelePresence endpoints (these values are exclusive of network overhead). These additions have been highlighted in bold in [Table 4-2](#). When the 20% network-overhead overprovisioning rule is applied, the additional bandwidth required to support this interoperability feature becomes about 1 Mbps for any TelePresence system, regardless of the number of segments, resolution-levels, and motion-handling capabilities configured for TelePresence primary video.

Table 4-2 Cisco TelePresence Software Version 1.3 Bandwidth Requirements (Including the Interoperability Feature)

Resolution	1080p	1080p	1080p	720p	720p	720p
Motion Handling	Best	Better	Good	Best	Better	Good
Video per Screen (kbps)	4000	3500	3000	2250	1500	1000
Audio per Microphone (kbps)	64	64	64	64	64	64
(5 fps) Auto Collaborate Video Channel (kbps)	500	500	500	500	500	500
Auto Collaborate Audio Channel (kbps)	64	64	64	64	64	64
Interoperability Video Channel (kbps)	768	768	768	768	768	768
Interoperability Audio Channel (kbps)	64	64	64	64	64	64
CTS-500/1000 Total Audio and Video (kbps)	5,460	4,960	4,460	3,710	2,960	2,460
CTS-3000/3200 Total Audio and Video (kbps)	13,588	12,088	10,588	8,338	6,088	4,588
CTS-500/1000 total bandwidth (Including Layer 2-Layer 4 overhead)	6.5 Mbps	6.0 Mbps	5.3 Mbps	4.4 Mbps	3.5 Mbps	3 Mbps
CTS-3000/3200 total bandwidth (Including Layer 2-Layer 4 overhead)	16.3 Mbps	14.5 Mbps	12.7 Mbps	10 Mbps	7.3 Mbps	5.5 Mbps

Subsequently, with the release of CTS software version 1.4, the auxiliary video/Auto Collaborate video channel was expanded to transmit 30 frame-per-second of video (from the previous level of 5 frames-per-second) by leveraging an optional, dedicated presentation codec. The 30 fps Auto Collaborate video channel requires 4 Mbps of video bandwidth (along with 64 kbps of audio bandwidth). In comparison, the 5 fps Auto Collaboration feature required only 500 kbps of video bandwidth (along with 64 kbps of audio bandwidth). This change has been highlighted in bold in [Table 4-3](#). Therefore, the net increase in bandwidth required to support this 30 fps Auto Collaborate feature is 3.5 Mbps (exclusive of network overhead). When the 20% network-overhead overprovisioning rule is applied, the additional bandwidth required to support this 30 fps Auto Collaborate feature becomes about 4.2 Mbps for any TelePresence system, regardless of the number of segments, resolution-levels, and motion-handling capabilities configured for TelePresence primary video.

Table 4-3 Cisco TelePresence Software Version 1.4 Bandwidth Requirements (Including the 30 fps Auto Collaborate Feature)

Resolution	1080p	1080p	1080p	720p	720p	720p
Motion Handling	Best	Better	Good	Best	Better	Good
Video per Screen (kbps)	4000	3500	3000	2250	1500	1000
Audio per Microphone (kbps)	64	64	64	64	64	64
(30 fps) Auto Collaborate Video Channel (kbps)	4000	4000	4000	4000	4000	4000
Audio Add-In channel (kbps)	64	64	64	64	64	64
Interoperability Video Channel (kbps)	768	768	768	768	768	768
Interoperability Audio Channel (kbps)	64	64	64	64	64	64
CTS-500/1000 Total Audio and Video (kbps)	8,960	8,460	7,960	7,210	6,460	5,960
CTS-3000/3200 Total Audio and Video (kbps)	17,088	15,588	14,088	11,838	9,588	8,088
CTS-500/1000 total bandwidth (Including Layer 2-Layer 4 overhead)	10.9 Mbps	10.2 Mbps	9.5 Mbps	8.6 Mbps	7.7 Mbps	7.2 Mbps
CTS-3000/3200 total bandwidth (Including Layer 2-Layer 4 overhead)	20.5 Mbps	18.7 Mbps	16.9 Mbps	14.2 Mbps	11.5 Mbps	9.7 Mbps

In conclusion, it bears repeating that as Cisco TelePresence software continues to evolve and add new feature support, the bandwidth requirements for TelePresence will correspondingly evolve and expand.

Burst Requirements

So far, we have discussed bandwidth in terms of bits per second (i.e., how much traffic is sent over a one second interval). However, when provisioning bandwidth and configuring queuing, shaping, and policing commands on routers and switches, burst must also be taken into account. Burst is defined as the amount of traffic (generally measured in bytes) transmitted per millisecond which exceeds the per-second average. For example, a CTS-3000 running at 1080p-Best at approximately 15 Mbps divides evenly into approximately 1,966 bytes per millisecond ($15 \text{ Mbps} \div 1,000 \text{ milliseconds}$).

Cisco TelePresence operates at 30 frames per second. This means that every 33ms a video frame is transmitted; we refer to this as a frame interval. Each frame consists of several thousand bytes of video payload, and therefore each frame interval consists of several dozen packets, with an average packet size of 1,100 bytes per packet. However, because video is variable in size (due to the variability of motion in the encoded video), the packets transmitted by the codec are not spaced evenly over each 33ms frame interval, but rather are transmitted in bursts measured in shorter intervals. Therefore, while the overall bandwidth (maximum) averages out to 15 Mbps over one second, when measured on a per millisecond basis the packet transmission rate is highly variable, and the number of bytes transmitted per millisecond for a 15 Mbps per second call bursts well above the 1,966 bytes per millisecond average. Therefore, adequate burst tolerance must be accommodated by all switch and router interfaces in the path (platform-specific recommendations are detailed in the subsequent design chapters).

TelePresence Latency Requirements

Cisco TelePresence has a network latency target of 150 ms; this target does not include codec processing time, but purely network flight time.

There may be scenarios, however, where this latency target may not always be possible to achieve, simply due to the laws of physics and the geographical distances involved. Therefore, TelePresence codecs have been designed to sustain high levels of call quality even up to 200 ms of latency. Beyond this threshold (which we refer to as ‘Latency Threshold 1’) a warning message appears on the screen indicating that network conditions may be affecting call quality. Nonetheless, the call continues. If network latency exceeds 400 ms (which we refer to as ‘Latency Threshold 2’) another warning message appears on the screen and the call quality steadily degrades as latency increases. Visually, the call quality is the same, but aurally the lagtime between one party speaking and the other party responding becomes unnaturally excessive. In the original release of the TelePresence codec, calls were self-terminated by the codec if network latency increased beyond 400 ms. However, due to some unique customer requirements, such as some customers looking at provisioning TelePresence calls over satellite circuits, this behavior changed for release 1.1 of the codec, in which the calls were no longer terminated if Latency Threshold 2 was exceeded. Nonetheless, should customers choose to provision TelePresence over such circuits, user expectations need to be adjusted accordingly.

Network latency time can be broken down further into fixed and variable components:

- Serialization (fixed)
- Propagation (fixed)
- Queuing (variable)

Serialization refers to the time it takes to convert a Layer 2 frame into Layer 1 electrical or optical pulses onto the transmission media. Therefore, serialization delay is fixed and is a function of the line rate (i.e., the clock speed of the link). For example, a 45 Mbps DS3 circuit would require 266 μ s to serialize a 1500 byte Ethernet frame onto the wire. At the circuit speeds required for TelePresence (generally speaking DS3 or higher), serialization delay is not a significant factor in the overall latency budget.

The most significant network factor in meeting the latency targets for TelePresence is propagation delay, which can account for over 90% of the network latency time budget. Propagation delay is also a fixed component and is a function of the physical distance that the signals have to travel between the originating endpoint and the receiving endpoint. The gating factor for propagation delay is the speed of light: 300,000 km/s or 186,000 miles per second. Roughly speaking, the speed of light in an optical fiber is slightly less than one third the speed of light in a vacuum. Thus, the propagation delay works out to be approximately 6.3 μ s per km or 8.2 μ s per mile.

Another point to keep in mind when calculating propagation delay is that optical fibers are not always physically placed over the shortest path between two geographic points, especially over transoceanic links. Due to installation convenience, circuits may be hundreds or thousands of kilometers longer than theoretically necessary.

Nonetheless, the network flight-time budget of 150 ms allows for nearly 24,000 km or 15,000 miles worth of propagation delay (which is approximately 60% of the earth’s circumference); the theoretical worst-case scenario (exactly half of the earth’s circumference) would require only 126 ms. Therefore, this latency target should be achievable for virtually any two locations on the planet, given relatively direct transmission paths. However, for some of the more extreme scenarios, user expectations may have to be set accordingly, as there is little a network administrator can do about increasing the speed of light.

Given the end-to-end latency targets and thresholds for TelePresence, the network administrator also must know how much of this budget is to be allocated to the service provider and how much to the enterprise. The general recommendation for this split is 80:20, with 80% of the latency budget allocated to the service provider (demarc-to-demarc) and 20% to the enterprise (codec-to-demarc on one side and

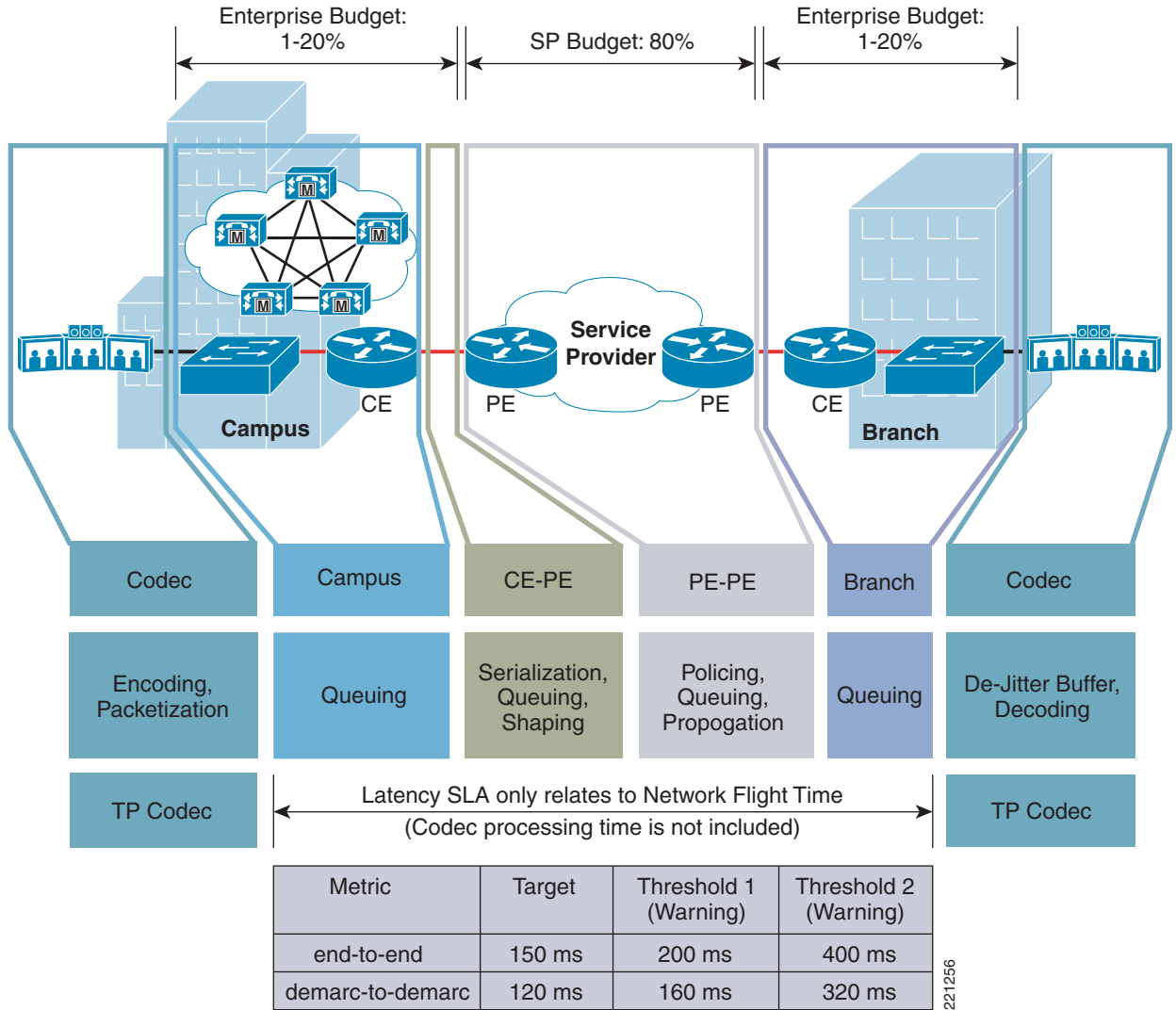
demarc-to-codec on the other). However, some enterprise networks may not require a full 20% of the latency budget and thus may reallocate their allowance to a 90:10 service provider-to-enterprise split, or whatever the case may be. The main point is that a fixed budget needs to be clearly apportioned to both the service provider and to the enterprise, such that the network administrators can design their networks accordingly. Given the target (150ms), threshold1 (200ms), and the service provider-enterprise split of 80:20 or 90:10, it is recommended that SPs engineer their network to meet the target, but base their SLA on threshold1. Threshold1 provides global coverage between any two sites on the planet and allows the SP to offer a 100% guarantee that their network (demarc-to-demarc) will never exceed 160ms (80% of threshold1).

Another point to bear in mind here is the additional latency introduced by multipoint resources. Latency is always measured from end-to-end (i.e., from codec1 to codec2). However, in a multipoint call the media between the two codecs traverses a Multipoint Switch. The multipoint switch itself introduces approximately 20ms of latency, and the path from codec1 to the MS and from the MS to codec2 may be greater than the path between codec1 and codec2 directly, depending on the physical location of the MS. Therefore, when engineering the network with respect to latency, one must calculate both scenarios for every TelePresence System deployed: one for the path between each system and every other system for point-to-point call, and a second for the path between each system, through the MS, to every other system.

The final TelePresence latency component to be considered is queuing delay, which is variable. Queuing delay is a function of whether a network node is congested and what the scheduling QoS policies are to resolve congestion events. Given that the latency target for TelePresence is very tight and, as has been shown, the majority of factors contributing to the latency budget are fixed, careful attention has to be given to queuing delay, as this is the only latency factor that is directly under the network administrator's control via QoS policies.

The latency targets, thresholds and service provider-to-enterprise splits are illustrated in [Figure 4-2](#).

Figure 4-2 Network Latency Target and Thresholds for Cisco TelePresence



TelePresence Jitter Requirements

Cisco TelePresence has a peak-to-peak jitter target of 10 ms. Jitter is defined as the variance in network latency. Thus, if the average latency is 100 ms and packets are arriving between 95 ms and 105 ms, the peak-to-peak jitter is defined as 10 ms. Measurements within the Cisco TelePresence codecs use peak-to-peak jitter.

Similar to the latency service level requirement, Cisco TelePresence codecs have built in thresholds for jitter to ensure a high quality user experience. Specifically, if peak-to-peak jitter exceeds 20 ms (which we call Jitter Threshold 1) for several seconds, then two things occur:

- A warning message appears at the bottom of the 65" plasma display indicating that the network is experiencing congestion and that call quality may be affected.
- The TelePresence codecs downgrade to a lower level of motion handling quality within the given resolution.

As previously mentioned, Cisco TelePresence codecs have three levels of motion handling quality within a given resolution, specifically 720p-Good, 720p-Better, and 720p-Best and 1080p-Good, 1080p-Better, and 1080p-Best. Therefore, for example, if a call at 1080p-Best would exceed Jitter Threshold 1 (20 ms) for several seconds, the codec would display the warning message in and would downgrade the motion handling quality to 1080p-Good. Similarly a call at 720p-Best would downgrade to 720-Good. Incidentally, downgraded calls do not automatically upgrade should network conditions improve, because this could cause a “flapping” effect where the call upgrades and then downgrades again, over and over.

A second jitter threshold (Jitter Threshold 2) is also programmed into the TelePresence codecs, such that if peak-to-peak jitter exceeds 40 ms for several seconds, then two things occur. The TelePresence codecs:

- Self-terminate the call.
- Display an error message on the 7975G IP Phone indicating that the call was terminated due to excessive network congestion.

Finally, as with latency, the jitter budget is proportioned between the service provider and enterprise networks. Unfortunately, unlike latency or packet loss, peak-to-peak jitter is not necessarily cumulative. Nonetheless, simply for the sake of setting a jitter target for each party, the recommended peak-to-peak jitter split is 50/50 between the service provider and enterprise, such that each group of network administrators can design their networks to a clear set of jitter targets and thresholds. Also like latency, this split may be negotiated differently between the service provider and enterprise to meet certain unique scenarios, such as satellite connections. Again, the main point is that a fixed jitter budget needs to be clearly apportioned to both the service provider and to the enterprise, such that the end-to-end target and thresholds are not exceeded.

It is recommended that SPs engineer their network to meet the target, but base their SLA on threshold1. Threshold1 provides global coverage between any two sites on the planet and allows the SP to offer a 100% guarantee that their network (demarc-to-demarc) will never exceed 10ms of jitter (50% of threshold1).

The TelePresence Jitter targets and thresholds are summarized in [Table 4-4](#).

Table 4-4 *TelePresence Jitter Targets, Thresholds, and Service Provide/Enterpriser Splits*

Metric	Target	Threshold 1 (Warning and Downgrade)	Threshold 2 (Call Drop)
End-to-end	10 ms	20 ms	40 ms
Service Provider	5 ms	10 ms ¹	20 ms

1. SP SLA should be based on Threshold 1.

TelePresence Loss Requirements

Cisco TelePresence is highly sensitive to packet loss, and as such has an end-to-end packet loss target of 0.05%.

It may be helpful to review a bit of background information to better understand why TelePresence is so extremely sensitive to packet loss. Specifically, let’s review how much information is actually needed to transmit a 1080p30 HD video image, which is the highest video transmission format used by Cisco TelePresence codecs. The first parameter (1080) refers to 1080 lines of horizontal resolution, which are matrixed with 1920 lines of vertical resolution (as per the 16:9 Widescreen Aspect Ratio used in High Definition video formatting), resulting in 2,073,600 pixels per screen. The second parameter, p, indicates a progressive scan, which means that every line of resolution is refreshed with each frame (as opposed to an interlaced scan, which would be indicated with an i and would mean that every other line is

refreshed with each frame). The third parameter 30 refers to the transmission rate of 30 frames per second. While video sampling techniques may vary, each pixel has approximately 3 Bytes of color and/or luminance information. When all of this information is factored together (2,073,600 pixels x 3 Bytes x 8 bits per Byte x 30 frames per second), it results in approximately 1.5 Gbps of information. This is illustrated in Figure 4-3.

Figure 4-3 1080p30 Information Breakdown



As shown earlier in this chapter, Cisco TelePresence codecs transmit at approximately 5 Mbps (max) per 1080p display, which translates to over 99% compression. Therefore, the overall effect of packet loss is proportionally magnified and dropping even one packet in 2000 (0.05% packet loss) becomes readily noticeable to end users.

Similar to the latency and jitter service level requirement, Cisco TelePresence codecs have built in thresholds for packet loss to ensure a high-quality user experience. Specifically, if packet loss exceeds 0.10% (or 1 in 1000 packets, which we call Loss Threshold 1) for several seconds, then two things occur:

- A warning message appears at the bottom of the on the 65" plasma display indicating that the network is experiencing congestion and that call quality may be affected.
- The TelePresence codecs downgrade to a lower level of motion handling quality within the given resolution.

As previously mentioned, Cisco TelePresence codecs have three levels of motion handling quality within a given resolution, specifically 720p-Good, 720p-Better, and 720p-Best and 1080p-Good, 1080p-Better, and 1080p-Best. Therefore, for example, if a call at 1080p-Best would exceed Loss Threshold 1 (0.10%) for several seconds, the codec would display the warning message and would downgrade the motion handling quality to 1080p-Good. Similarly a call at 720p-Best would downgrade to 720p-Good in the same scenario. Incidentally, downgraded calls do not automatically upgrade should network conditions improve, because this could cause a "flapping" effect where the call upgrades and then downgrades again, over and over.

A second packet loss threshold (Loss Threshold 2) is also programmed into the TelePresence codecs, such that if packet loss exceeds 0.20% (or 1 in 500 packets) for several seconds, then two things occur. The TelePresence codecs:

- Self-terminate the call.
- Display an error message on the 7975G IP Phone indicating that the call was terminated due to excessive network congestion.

Finally, as with previously defined service level requirements, the loss budget is proportioned between the service provider and enterprise networks. The recommend split is 50/50 between the service provider and enterprise, such that each group of network administrators can design their networks to a clear set of packet loss targets and thresholds. Of course, This split may be negotiated differently between the service provider and enterprise to meet certain unique scenarios, such as satellite connections. Again, the main point is that a fixed packet loss budget needs to be clearly apportioned to both the service provider and to the enterprise, such that the end-to-end target and thresholds are not exceeded.

It is recommended that SPs engineer their network to meet the target, but base their SLA on threshold1. Threshold1 provides global coverage between any two sites on the planet and allows the SP to offer a 100% guarantee that their network (demarc-to-demarc) will never exceed .05% loss (50% of threshold1).

The TelePresence packet loss targets and thresholds are summarized in [Table 4-5](#).

Table 4-5 *TelePresence Jitter Targets, Thresholds, and Service Provider/Enterprise Splits*

Metric	Target	Threshold 1 Warning and Downgrade)	Threshold 2 (Call Drop)
End-to-end	0.05% (1 in 2000)	0.10% (1 in 1000)	0.20 (1 in 500)
Service Provider	.025%	.05% ¹	.10%

1. SP SLA should be based on Threshold 1.

Tactical QoS Design Best Practices for TelePresence

Once the service level requirements of TelePresence are defined, then the network administrator can proceed to the next step of the QoS deployment cycle (illustrated in [Figure 4-1](#)) of designing the actual policies.

A couple of tactical QoS best practices design principles bear mentioning at this point, as these serve to improve the efficiency and scope of your QoS designs. The first principle is to always deploy QoS in hardware, rather than software, whenever a choice exists. Cisco Catalyst switches perform QoS operations in hardware Application Specific Integrated Circuits (ASICs) and as such have zero CPU impact; Cisco IOS routers, on the other hand, perform QoS operations in software, resulting in a marginal CPU impact, the degree of which depends on the platform, the policies, the link speeds, and the traffic flows involved. So, whenever supported, QoS policies like classification, marking/remarking, and/or policing can all be performed at line rates with zero CPU impact in Catalyst switches (as opposed to IOS routers), which makes the overall QoS design more efficient. A practical example of how this principle is applied is as follows: while all nodes in the network path must implement queuing policies, classification policies should be implemented in Cisco Catalyst hardware as close to the source of the traffic as possible (e.g., on the access edge switch to which the TelePresence System is attached), rather than waiting until the traffic hits the WAN router to be classified.

Another best practice principle to keep in mind is to follow industry standards whenever possible, as this extends the effectiveness of your QoS policies beyond your direct administrative control. For example, if you mark a realtime application, such as VoIP, to the industry standard recommendation as defined in RFC 3246 (An Expedited Forwarding Per-Hop Behavior), then you will no doubt provision it with strict priority servicing at every node within your enterprise network. Additionally, if you handoff to a service provider following this same industry standard, they will similarly provision traffic marked Expedited Forwarding (EF - or DSCP 46) in a strict priority manner. Therefore, even though you do not have direct administrative control of the QoS policies within the service provider's cloud, you have extended the influence of your QoS design to include your service provider's cloud, simply by following the industry

standard recommendations. Therefore, in line with this principle, it would be beneficial to briefly consider some of the relevant industry standards to QoS design, particularly as these relate to TelePresence.

Relevant Industry Standards and Recommendations

Let's briefly review some of the relevant DiffServ standards and recommendations and see how these relate to TelePresence QoS design.



Note

Although Cisco TelePresence requires Cisco CallManager (CCM) 5.1 (or higher) for call processing, and CCM 5.x supports Resource Reservation Protocol (RSVP) for Call Admission Control, the initial phase of the TelePresence solution does not require leveraging RSVP functionality (RSVP remains optional during this phase); therefore, the discussion in this paper focuses on DiffServ QoS designs and standards for Cisco TelePresence (not IntServ/RSVP).

RFC 2474 Class Selector Code Points

This standard defines the use of 6 bits in the IPv4 and IPv4 Type of Service (ToS) byte, termed Differentiated Services Code Points (DSCP). Additionally, this standard introduces Class Selector codepoints to provide backwards compatibility for legacy (RFC 791) IP Precedence bits.

RFC 2597 Assured Forwarding Per-Hop Behavior Group

This standard defines the Per-Hop Behavior of the Assured Forwarding (AF) classes. Four AF classes are defined: AF1, AF2, AF3, and AF4. Additionally, each class has three states of increasing Drop Preference assigned within it, corresponding to three traffic states: conforming (analogous to a green traffic light signal), exceeding (analogous to a yellow traffic light signal), and violating (analogous to a red traffic light signal). For example, conforming AF1 traffic would be marked to AF11 (the second 1 representing the lowest Drop Preference setting), exceeding traffic would have its Drop Preference increased to AF12, and violating traffic would have its Drop Preference increased further to AF13. When such traffic enters a node experiencing congestion, AF13 traffic is more aggressively dropped than AF12 traffic, which in turn is more aggressively dropped than AF11 traffic.

RFC 3246 An Expedited Forwarding Per-Hop Behavior

This standard defines an Expedited Forwarding (EF) Per-Hop Behavior for realtime applications. When traffic marked EF enters a node experiencing congestion, it receives strict priority behavior.

RFC 3662 A Lower Effort Per-Domain Behavior for Differentiated Services

This informational RFC defines a less than Best Effort service for undesired applications and specifies that such applications should be marked to Class Selector 1 (CS1).

Cisco's QoS Baseline

While the IETF RFC standards provided a consistent set of per-hop behaviors for applications marked to specific DSCP values, they never specified which application should be marked to which DiffServ Codepoint value. Much confusion and disagreements over matching applications with standards-defined

codepoints led Cisco in 2002 to put forward a standards-based marking recommendation in their strategic architectural QoS Baseline document. Eleven different application classes that could exist within the enterprise were examined and extensively profiled, and then matched to their optimal RFC-defined Per-Hop Behaviors (PHBs). The application-specific marking recommendations from Cisco's QoS Baseline of 2002 are summarized in [Figure 4-4](#).

Figure 4-4 Cisco's QoS Baseline Marking Recommendations

Application	L3 Classification		IETF
	PHB	DSCP	RFC
Routing	CS6	48	RFC 2474
Voice	EF	46	RFC 3246
Interactive Video	AF41	34	RFC 2597
Streaming Video	CS4	32	RFC 2474
Mission-Critical Data	AF31	26	RFC 2597
Call Signaling	CS3	24	RFC 2474
Transactional Data	AF21	18	RFC 2597
Network Management	CS2	16	RFC 2474
Bulk Data	AF11	10	RFC 2597
Best Effort	0	0	RFC 2474
Scavenger	CS1	8	RFC 2474

220199

The adoption of Cisco's QoS Baseline was a great step forward in QoS consistency, not only within Cisco, but also within the industry in general.

RFC 4594 Configuration Guidelines for DiffServ Classes

More than four years after Cisco put forward its QoS Baseline document, RFC 4594 was formally accepted as an informational RFC (in August 2006).

Before getting into the specifics of RFC 4594, it is important to comment on the difference between the IETF RFC categories of informational and standard. An informational RFC is an industry recommended best practice, while a standard RFC is an industry requirement. Therefore RFC 4594 is a set of formal DiffServ QoS configuration best practices, not a requisite standard.

RFC 4594 puts forward twelve application classes and matches these to RFC-defined Per-Hop Behaviors (PHBs). These application classes and recommended PHBs are summarized in [Figure 4-5](#).

Figure 4-5 RFC 4594 Marking Recommendations

Application	L3 Classification		IETF
	PHB	DSCP	RFC
Network Control	CS6	48	RFC 2474
VoIP Telephony	EF	46	RFC 3246
Call Signaling	CS5	40	RFC 2474
Multimedia Conferencing	AF41	34	RFC 2597
Real-Time Interactive	CS4	32	RFC 2474
Multimedia Streaming	AF31	26	RFC 2597
Broadcast Video	CS3	24	RFC 2474
Low-Latency Data	AF21	18	RFC 2597
OAM	CS2	16	RFC 2474
High-Throughput Data	AF11	10	RFC 2597
Best Effort	DF	0	RFC 2474
Low-Priority Data	CS1	8	RFC 3662

220200

It is fairly obvious that there are more than a few similarities between Cisco's QoS Baseline and RFC 4594, as there should be, since RFC 4594 is essentially an industry-accepted evolution of Cisco's QoS Baseline. However, there are some differences that merit attention.

The first set of differences are minor, as they involve mainly nomenclature. Some of the application classes from the QoS Baseline have had their names changed in RFC 4594. These changes in nomenclature are summarized in [Table 4-6](#).

Table 4-6 Nomenclature Changes from Cisco QoS Baseline to RFC 4594

Cisco QoS Baseline Class Names	RFC 4594 Class Names
Routing	Network Control
Voice	VoIP Telephony
Interactive Video	Multimedia Conferencing
Streaming Video	Multimedia Streaming
Transactional Data	Low-Latency Data
Network Management	Operations/Administration/Management (OAM)
Bulk Data	High-Throughput Data
Scavenger	Low-Priority Data

The remaining changes are more significant. These include one application class deletion, two marking changes, and two new application class additions. Specifically:

- The QoS Baseline Locally-Defined Mission-Critical Data class has been deleted from RFC 4594.
- The QoS Baseline marking recommendation of CS4 for Streaming Video has been changed in RFC 4594 to mark Multimedia Streaming to AF31.

- The QoS Baseline marking recommendation of CS3 for Call Signaling has been changed in RFC 4594 to mark Call Signaling to CS5.
- A new video class has been added to RFC 4594: Real-Time Interactive, which is to be marked CS4. This was done to differentiate between lower-grade desktop video telephony (referred to as Multimedia Conferencing) and higher-grade videoconferencing and TelePresence. Multimedia Conferencing uses the AF4 class and is subject to markdown policies, while TelePresence uses the CS4 class and is not subject to markdown.
- A second new video class has been added to RFC 4594: Broadcast video, which is to be marked CS3. This was done to differentiate between lower-grade desktop video streaming (referred to as Multimedia Streaming) and higher-grade Broadcast Video applications. Multimedia Streaming uses the AF3 class and is subject to markdown policies, while Broadcast Video uses the CS3 class and is not subject to markdown.

The most significant of the differences between Cisco's QoS Baseline and RFC 4594 is the RFC 4594 recommendation to mark Call Signaling to CS5. Cisco has just completed a lengthy and expensive marking migration for Call Signaling from AF31 to CS3 (as per the original QoS Baseline of 2002), and as such, there are no plans to embark on another marking migration in the near future. It is important to remember that RFC 4594 is an informational RFC (i.e., an industry best-practice) and not a standard. Therefore, lacking a compelling business case at the time of writing, Cisco plans to continue marking Call Signaling as CS3 until future business requirements arise that necessitate another marking migration.

Therefore, for the remainder of this document, RFC 4594 marking values are used throughout, with the one exception of swapping Call-Signaling marking (to CS3) and Broadcast Video (to CS5). These marking values are summarized in [Figure 4-6](#).

Figure 4-6 Cisco-Modified RFC4594 Marking Values (Call-Signaling is Swapped with Broadcast Video)

Application	L3 Classification		IETF
	PHB	DSCP	RFC
Network Control	CS6	48	RFC 2474
VoIP Telephony	EF	46	RFC 3246
Broadcast Video	CS5	40	RFC 2474
Multimedia Conferencing	AF41	34	RFC 2597
Real-Time Interactive	CS4	32	RFC 2474
Multimedia Streaming	AF31	26	RFC 2597
Call Signaling	CS3	24	RFC 2474
Low-Latency Data	AF21	18	RFC 2597
OAM	CS2	16	RFC 2474
High-Troughput Data	AF11	10	RFC 2597
Best Effort	DF	0	RFC 2474
Low-Priority Data	CS1	8	RFC 3662

221258

Classifying TelePresence

One of the first questions to be answered relating to TelePresence QoS design is: should TelePresence be assigned to a dedicated class or should it be assigned to the same class as existing Videoconferencing/Video Telephony? The answer to this question directly relates to whether TelePresence has the same service-level requirements as these other two interactive video applications or whether it has unique service level requirements. [Table 4-7](#) summarizes the service level requirements of both generic Videoconferencing applications and TelePresence.

Table 4-7 Service Level Requirements of Generic Video-Conferencing and TelePresence

Service Level Parameter (Target Values)	(Generic) Videoconferencing/Video Telephony	Cisco TelePresence
Bandwidth	384 kbps or 768 kbps + network overhead	1.5 Mbps to 12.6 Mbps + network overhead
Latency	400-450 ms latency	150 ms latency
Jitter	30-50 ms peak-to-peak jitter	10 ms peak-to-peak jitter
Loss	1% random packet loss	0.05% random packet loss

From [Table 4-7](#) it becomes apparent that TelePresence has unique (and higher/tighter) service level requirements than do generic Videoconferencing/Video Telephony applications; therefore, TelePresence requires a dedicated class along with a dedicated classification marking value.

Videoconferencing/Video Telephony applications have traditionally been marked to (RFC 2597) Assured Forwarding Class 4, which is the recommendation from both the Cisco QoS Baseline as well as RFC 4594. However, the Assured Forwarding (AF) Per-Hop Behavior (PHB) includes policing (to conforming, exceeding, and violating traffic rates), as well as correspondingly increasing the Drop Preferences (to Drop Preference 1, 2, and 3 respectively), and ultimately dropping traffic according to the Drop Preference markings. TelePresence traffic has a very low tolerance to drops (0.05%) and therefore would not be appropriately serviced by an AF PHB.

Because of the low-latency and jitter service-level requirements of TelePresence, it may seem attractive to assign it an (RFC 3246) Expedite Forwarding (EF) Per-Hop Behavior; after all, there is nothing in RFC 3246 that dictates that only VoIP can be assigned to this PHB. However, it is important to recognize that VoIP behaves considerably differently than video. As previously mentioned, VoIP has constant packet sizes and packet rates, whereas video packet sizes vary and video packet rates also vary in a random and bursty manner. Thus, if both video and voice were assigned to the same marking value and class, (bursty) video could easily interfere with (well-behaved) voice. Therefore, for both operational and capacity planning purposes, it is recommended not to mark both voice and video to EF. This recommendation is reflected in both the Cisco QoS Baseline as well as RFC 4594.

What then should TelePresence be marked to? The best formal guidance is provided in RFC 4594, where a distinction is made between a Multimedia Conferencing (i.e., generic Videoconferencing/Video Telephony) service class and a Real-Time Interactive service class. The Real-Time Interactive service class is intended for inelastic video flows, such as TelePresence. The recommended marking for this Real-Time Interactive service class, and thus **the recommended marking for TelePresence is Class Selector 4 (CS4)**.

Policing TelePresence

In general, policing TelePresence traffic should be avoided whenever possible, although some exceptions exist.

As previously mentioned, TelePresence is highly sensitive to drops (with a 0.05% packet loss target); therefore policing TelePresence traffic rates with either a Single Rate Three Color Marker (as defined in RFC 2697) or a Two Rate Three Color Marker (as defined in RFC 2698) could be extremely detrimental to TelePresence flows and ultimately ruin the high-level of user experience that this application is intended to deliver.

However, there are three places where TelePresence traffic may be legitimately policed over the network.

The first automatically occurs if TelePresence is assigned to a Low-Latency Queue (LLQ) within Cisco IOS routers at the WAN or VPN edge. This is because any traffic assigned to a LLQ is automatically policed by an implicit policer set to the exact value as the LLQ rate. For example, if TelePresence is assigned a LLQ of 15 Mbps, it is also implicitly policed by the LLQ algorithm to exactly 15 Mbps; any excess TelePresence traffic is dropped.



Note

The implicit policer within the LLQ feature is only active when LLQ is active. In other words, since queuing only engages when there is congestion, LLQ never engages unless the link is physically congested or a (hierarchical QoS) shaper forces LLQ to engage prior to physical link congestion. Similarly, the implicit policer of LLQ never engages unless there is physical congestion on the link or a (hierarchical QoS) shaper forces it to engage prior to physical link congestion. Put another way, when the physical link is un-congested and/or a hierarchical QoS shaper is inactive, neither LLQ nor the implicit policer of LLQ is active.

The second most common place that TelePresence is likely to be policed in the network is at the service provider's provider edge (PE) routers, in the ingress direction. Service providers need to police traffic classes, especially realtime traffic classes, to enforce service contracts and prevent possible oversubscription on their networks and thus ensure service level agreements.

The third place (and optional) place, where policing TelePresence may prove beneficial in the network is at the campus access edge. Administrators can deploy access-edge policers for security purposes to mitigate the damage caused by the potential abuse of trusted switch ports. Since TelePresence endpoints can mark TelePresence flows to the recommended 802.1Q/p CoS value (CoS 4) and DSCP codepoint value (CS4), the network administrator may choose to trust the CoS or DSCP values received from these ports. However, if a disgruntled employee gains physical access to the TelePresence switch ports, they may send whatever traffic they choose to over these ports and their flows are trusted over the network. Such rogue traffic flows may hijack voice or video queues and easily ruin call or video quality over the QoS-provisioned network infrastructure. Therefore, the administrator may choose to limit the scope of damage that such network abuse may present by configuring access-edge policers on TelePresence switch ports to remark (to Scavenger: DSCP CS1) or drop out-of-profile traffic originating on these ports (e.g., CS4 traffic exceeding 15 Mbps). Supporting this approach, RFC 4594 recommends edge policing the Real-Time Interactive service class via a single-rate policer.

Queuing TelePresence

To achieve the high-levels of service required by the Cisco TelePresence Experience, queuing must be enabled on every node along the path to provide service guarantees, regardless of how infrequently congestion may occur on certain nodes (i.e., congestion can and does occur even on very high-bandwidth mediums). If queuing is not properly configured on every node, the Cisco TelePresence eXperience (CTX) cannot be guaranteed.

RFC 4594 specifies the minimum queuing requirement of the Real-Time Interactive service class to be a rate-based queue (i.e., a queue that has a guaranteed minimum bandwidth rate). However, RFC 4594 makes an allowance that while **the PHB for Real-Time Interactive service class** should be configured to provide high bandwidth assurance, it **may be configured as a second EF PHB** that uses relaxed performance parameters, a rate scheduler, and a CS4 DSCP value.

This means that, for example, TelePresence, which has been assigned to this Real-Time Interactive service class, can be queued with either a guaranteed rate non-priority queue (such as a Cisco IOS Class-Based Weighted Fair Queue-CBWFQ) or a guaranteed-rate strict priority queue (such as a Cisco IOS Low-Latency Queue-LLQ); in either case, TelePresence is to be marked as Class Selector 4 (and not EF).

Therefore, since RFC 4594 allows for the Real-Time Interactive service-class to be given a second EF PHB and because of the low latency, low jitter, and low loss requirements of TelePresence, **it is recommended to place TelePresence in a strict-priority queue**, such as a Cisco IOS LLQ or a Cisco Catalyst hardware priority queue whenever possible.

However, an additional provisioning consideration must be taken into account when provisioning TelePresence with a second EF PHB, which relates to the amount of bandwidth of a given link that should be assigned for strict priority queuing. The well-established and widely-deployed Cisco best-practice recommendation is to limit the amount of strict priority queuing configured on an interface to no more than one-third of the link's capacity. This has commonly been referred to as the 33% LLQ Rule.

The rationale behind this rule is that if you assign too much traffic for strict priority queuing, then the overall effect is a dampening of QoS functionality for non-realtime applications. Remember, the goal of convergence is to enable voice, video, and data to transparently co-exist on a single network. When realtime applications such as voice and/or TelePresence dominate a link (especially a WAN/VPN link), then data applications fluctuate significantly in their response times when TelePresence calls are present versus when they are absent, thus destroying the transparency of the converged network.

For example, consider a (45 Mbps) DS3 link configured to support 2 separate CTS-3000 calls, both configured to transmit at full 1080p-Best resolution. Each such call requires 15 Mbps of realtime traffic. Prior to TelePresence calls being placed, data applications have access to 100% of the bandwidth (to simplify the example, we are assuming there are no other realtime applications, such as VoIP, on this link). However, once these TelePresence calls are established, all data applications would suddenly be contending for less than 33% of the link. TCP windowing would take effect and many data applications will hang, time-out, or become stuck in a non-responsive state, which usually translates into users calling the IT help desk complaining about the network (which happens to be functioning properly, albeit in a poorly-configured manner).

To obviate such scenarios, Cisco Technical Marketing has done extensive testing and has found that a significant decrease in data application response times occurs when realtime traffic exceeds one-third of link bandwidth capacity. Extensive testing and customer deployments have shown that a general best queuing practice is to limit the amount of strict priority queuing to 33% of link bandwidth capacity. This strict priority queuing rule is a conservative and safe design ratio for merging realtime applications with data applications.

**Note**

As Cisco IOS software allows the abstraction (and thus configuration) of multiple strict priority LLQs, in such a multiple LLQ context, this design principle would apply to the sum of all LLQs to be within one-third of link capacity.

It is vitally important, however, to understand that this strict priority queuing rule is simply a best practice design recommendation and is not a mandate. There may be cases where specific business objectives cannot be met while holding to this recommendation. In such cases, enterprises must provision according to their detailed requirements and constraints. However, it is important to recognize the tradeoffs involved with over-provisioning strict priority traffic and its negative performance impact on non-realtime-application response times. It is also worth noting that the 33% rule only applies for converged networks. In cases where customers choose to deploy dedicated WAN circuits for their TelePresence traffic, the 33% rule does not apply since TelePresence (and perhaps some nominal amount of management and signaling traffic) is the only traffic on the circuit. In these cases, customers are free to use up to 98% of the link capacity for TelePresence (reserving 2% for routing protocols, network management traffic such as SSH and SNMP, and signaling).

Shaping TelePresence?

It is recommended to avoid shaping TelePresence flows unless absolutely necessary. This is because of the QoS objective of shapers themselves. Specifically, the role of shapers is to delay traffic bursts above a certain rate and to smooth out flows to fall within contracted rates. Sometimes this is done to ensure traffic rates are within a carrier's Committed Information Rate (CIR); other times shaping is performed to protect other data classes from a bursty class.

Shapers temporarily buffer traffic bursts above a given rate and as such introduce variable delay (jitter) as well as absolute delay. Since TelePresence is so sensitive to delay (150 ms) and especially jitter (10 ms), it is recommended not to shape TelePresence flows.

If the objective of the shaper was to meet a carrier's CIRs, this can be achieved by properly provisioning the adequate bandwidth and burst allowances on the circuit.

If the objective of the shaper was to protect other traffic classes from TelePresence bursts, then a better approach would be to explicitly protect each class with a guaranteed minimum bandwidth rate (such as a Cisco IOS CBWFQ).

In either case, a shaper would be a sub-optimal tool to meet the desired objective and would cause quality issues on the TelePresence flows and therefore would not be recommended.

The TelePresence traffic queue (whether you choose to place it in a CBWFQ or a second strict priority LLQ) must be provisioned with the proper mean rate (bits per second) and burst allowance (burst bytes exceeding the mean).

Compressed RTP (cRTP) with TelePresence

It is recommended to not enable cRTP for TelePresence. This is because of the large CPU impact of IP RTP Header Compression and the negligible returns in bandwidth savings it entails at TelePresence circuit speeds.

TelePresence, like VoIP, is encapsulated by IP, UDP, and RTP headers and these headers, when combined, account for 40 bytes per packet (at Layer 3). To enhance bandwidth efficiency, compression tools, like IP RTP Header Compression (cRTP) can reduce this overhead from 40 bytes to 2-5 bytes per packet.

However, it is important to recognize that cRTP is the most computationally-intensive QoS operation in the Cisco IOS toolset. Furthermore, it is only recommended on slow-speed links, usually 768 kbps or less, as it is at these speeds that the bandwidth gain offsets the increased CPU cost of the operation and is only useful for RTP-based applications that have a small amount of payload per packet. On high-speeds links and applications like TelePresence in which the payload of each packet averages 1100 bytes, cRTP offers no benefit and only results in sending the routers CPU through the ceiling. **Therefore, it is recommended to not enable cRTP on links carrying TelePresence.**

Link Fragmentation and Interleaving (LFI) with TelePresence

Like cRTP, LFI is only useful on slow-speed links (usually 768 kbps or less) and is used to fragment larger data packets into smaller chunks and interleave voice in between them to reduce the serialization and queuing delays for VoIP applications. Since TelePresence packets average 1100 bytes payload per packet, LFI would want to fragment them. This introduces unwanted jitter and out-of-order and late packets into the TelePresence stream. On high-speed links the serialization delay for large packets is inconsequential to VoIP and thus LFI offers no benefit and only results in sending the routers CPU through the ceiling. **Therefore, it is recommended to not implement LFI on links carrying TelePresence.**

GRE/IPSec Tunnels with TelePresence

Tunneling TelePresence traffic over GRE/IPSec tunnels is supported. The Cisco TelePresence codecs are designed to limit their packets to a maximum of 1200 bytes to leave enough room for GRE/IPSec encapsulation overhead to avoid having the TelePresence traffic fragmented for exceeding the Maximum Transmission Unit (MTU) of any link in the path.

Place in the Network TelePresence QoS Design

At this point, the strategic QoS business objectives for TelePresence have been defined, the service level-requirements of TelePresence have been specified, and the tactical QoS design approach has been sketched via the best practice principles and recommendations reviewed in the previous section. What remains is to flesh out these sketches into detailed Place-in-the-Network (PIN) platform-specific designs.

As the Cisco TelePresence solution evolves, it will become more complex and touch more Places-in-the-Network. The first deployment model to receive Cisco Verified Design (CVD) certification is the Intra-Enterprise, Point-to-Point Deployment Model (as described in [Chapter 3, “TelePresence Network Deployment Models”](#)). Such deployments will directly impact enterprise campus, branch, and WAN/MAN PINs, as well as service provider edge and core networks.

An addition to the Intra-Enterprise Deployment Model came with the release of the Cisco TelePresence Multipoint Solution, based on the Cisco TelePresence Multipoint Switch (CTMS) product offering. This addition may require an additional PIN, namely the enterprise and/or service provider data center, as these are often the locations where multipoint resources are hosted. However, note that while many customers are beginning to deploy multipoint resources, the addition of multipoint resources within the Intra-Enterprise Deployment Model has not yet received CVD certification.

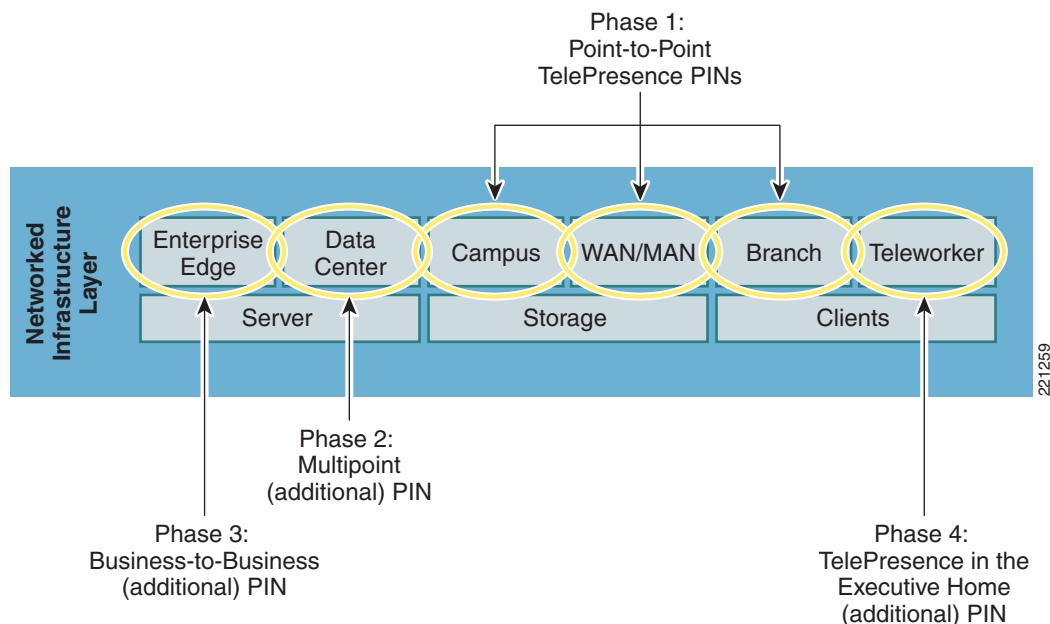
The next phase of TelePresence deployments will begin with the release of the Business-to-Business TelePresence solution, enabling enterprises to move to a Inter-Enterprise Deployment Model (as described in [Chapter 3, “TelePresence Network Deployment Models”](#)). These Inter-Enterprise deployments may be Point-to-Point or Multipoint. With this additional functionality, a new enterprise

PIN, the enterprise edge, will require design modifications. Additionally, service providers will need to develop shared services domains to provide the necessary connectivity, security, and QoS services required to enable this solution. Early Field Trials (EFT) of B2B services have begun. However, the Inter-Enterprise Deployment Model has not yet received CVD certification.

Finally, TelePresence systems are already emanating considerable executive-perk appeal, especially CTS-1000 systems that are designed for an executive's office. Already some executives are deploying TelePresence systems within their homes, taking advantage of very high-speed residential internet access options, like fiber optics to the home. Therefore, an inevitable fourth phase of TelePresence deployments will undoubtedly include the executive teleworker PIN. Early Field Trials (EFT) of TelePresence systems deployed in executive homes has begun. However, the Executive-Class Teleworker Deployment Model has not yet received CVD certification.

The relevant enterprise PINs for the above deployment models, based on the Service Oriented Network Architecture (SONA), specifically the Networked Infrastructure Layer, are illustrated in Figure 4-7.

Figure 4-7 SONA Networked Infrastructure Layer—Places in the Network (PINs) for Phases 1-4 TelePresence Deployments



The following chapters discuss and detail QoS designs for deploying TelePresence in each of these enterprise PINs. Information is provided on components pending CVD certification to allow customers to plan their network designs and deployment strategies accordingly. However, where detailed CVD design guidance is not yet available, note that the information provided is subject to change pending CVD certification.



CHAPTER 5

Campus QoS Design for TelePresence

Overview

The campus is the primary Place-in-the-Network (PIN) where TelePresence endpoints connect to the network infrastructure. Specifically, the 10/100/1000 NIC on the TelePresence primary codec connects—typically via an Intermediate Distribution Frame (IDF)—to the campus access layer edge switch port. It is at this switch port that the initial QoS polices required to support TelePresence are enabled. Additional QoS policies are also required on all campus inter-switch links, such as uplinks/downlinks between the access and distribution layers, as well as uplinks/downlinks from the distribution-to-core layers, and all core-layer links. Let's consider each of these port-specific QoS requirements.

Access Edge Switch Port QoS Considerations

The first QoS operation that needs to be performed is to define the trust boundary. The trust boundary is the point in the network at which 802.1Q/p CoS markings and/or IP DSCP markings are accepted or overridden by the network.

At the access-layer, the network administrator can enable the infrastructure to:

- Trust the endpoints (CoS and/or DSCP)
- Not trust the endpoints and manually re-mark TelePresence traffic using administratively-defined policies within the access-edge switch
- Conditionally trust the endpoints (trust is extended only after a successful CDP negotiation)

In Phase 1 deployments of Cisco TelePresence (Intra-Enterprise Point-to-Point Deployment Models, as discussed in [Chapter 3, “TelePresence Network Deployment Models”](#)) it is recommended to have a dedicated Communications Manager (CUCM) to support TelePresence. By default, CUCM marks any and all video traffic (including TelePresence) to AF41. It is recommended that this parameter be modified to mark video (i.e., TelePresence only, in this dedicated CUCM context) to CS4.



Note

The reason behind this recommendation is that CUCM does not (yet) have the ability to distinguish between different types of video. Therefore CUCM by default marks both generic Videoconferencing/Video Telephony (from applications like Cisco Unified Video Advantage, for example) as well as TelePresence to AF41.

If a dedicated CUCM is being used for managing TelePresence endpoints, and it has been configured to mark video (i.e., TelePresence) traffic to DSCP CS4, then the TelePresence primary codec marks all TelePresence call traffic (both video and audio) to CS4, but Call-Signaling traffic to CS3. The Cisco

7975G IP Phone similarly marks DSCP values correctly, marking VoIP traffic to EF and Call-Signaling to CS3. Therefore, the switch port connecting to the TelePresence primary codec can be configured to trust DSCP.

Alternatively, the access switch ports can be set to trust CoS, as both the Cisco 7975G IP Phone and the TelePresence primary codec are assigned to the Voice VLAN (VVLAN) and tag their traffic with 802.1Q/p CoS values. The 7975G IP Phone marks VoIP traffic to CoS 5 and Call-Signaling traffic to CoS 3. The Cisco TelePresence codec marks TelePresence traffic (both video and audio) to CoS 4 and Call-Signaling traffic to CoS 3.

However, if the switch port is configured to trust CoS, then it generates an internal DSCP value for all traffic flows via the CoS-to-DSCP map. Only one change is recommended to be made to the default CoS-to-DSCP map, which is to map CoS 5 to EF (46) instead of leaving the default mapping of CoS 5 to CS5 (40). The recommended CoS-to-DSCP map for access-switches connecting to Cisco TelePresence primary codecs is illustrated in [Table 5-1](#).

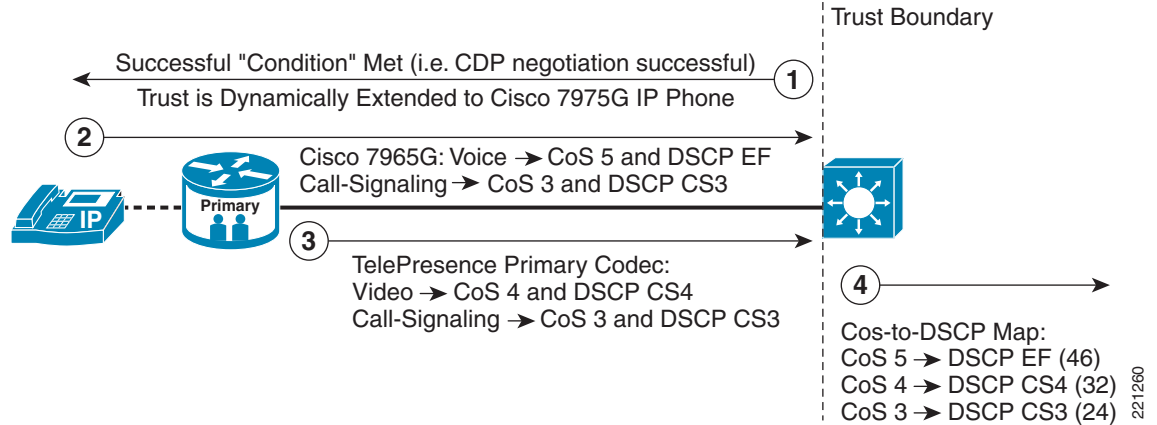
Table 5-1 Recommended Global CoS-to-DSCP Mapping for TelePresence Campus Switches

CoS Value	DSCP Value	PHB	Application
7	56	-	Network Control
6	48	CS6	Internetwork Control
5	46	EF	Voice
4	32	CS4	TelePresence
3	24	CS3	Call-Signaling
2	16	CS2	Management
1	8	CS1	Scavenger
0	0	DF	Default Forwarding/Best Effort

Finally, the access switch may be set to conditionally trust the TelePresence endpoint. This is because Cisco IP Telephony devices, including the Cisco Unified 7979G IP phone that is an intrinsic part of the TelePresence endpoint system, have the ability to identify themselves, via Cisco Discovery Protocol (CDP) to the network infrastructure. Upon a successful CDP negotiation/identification, the network infrastructure dynamically extends trust to the endpoints, which include both the Cisco Unified 7975G IP phone and the TelePresence primary codec. The primary functionality that conditional trust brings is to allow for user-mobility within the IP Telephony-enabled enterprise (users can add/move/change where their IP Phones are connected and the network automatically adapts without requiring an administrator to manually change switch port trust policies). This user-mobility is not a crucial functionality to support TelePresence, since TelePresence units are rarely moved around (due to sheer size). Nonetheless, this conditional trust functionality is supported by TelePresence codecs and adds a minor element of security in the event that the TelePresence codec is physically disconnected from the wall network jack by an unknowing and/or disgruntled individual, who then connects some other device (such as their laptop) to this trusted switch port. In this case, by using a conditional trust policy, the abuser's traffic would no longer be trusted.

The operation of conditional trust policies, as well as endpoint CoS markings and the CoS-to-DSCP mappings of the access-edge switch for TelePresence scenarios, is illustrated in [Figure 5-1](#).

Figure 5-1 Conditional Trust, CoS Markings, and Mappings for TelePresence



Note that if trust CoS is used (as opposed to conditional trust), steps 2, 3, and 4 still apply. The only difference is that the switch would skip step 1; the port would always be trusted regardless of CDP.

An optional recommendation for the access-edge switch port connecting to a TelePresence primary codec is to configure a policer to prevent network abuse in case of a compromise of this trusted port. Similar to the example previously given, this recommendation is to prevent an unknowing and/or disgruntled individual that gains physical access to the TelePresence switch port and decides to send rogue traffic over the network that can hijack voice or video queues and easily ruin call or video quality. Therefore, the administrator may choose to limit the scope of damage that such network abuse may present by configuring access-edge policers on TelePresence switch ports to drop (or remark to Scavenger - CS1) out-of-profile traffic originating on these ports. This is not only a Cisco recommended best practice, but is also reflected in RFC 4594 which recommends edge policing the Real-Time Interactive service class via a single-rate policer.

If such a policer is configured, it is recommended to use Per-Port/Per-VLAN policers, whenever supported. In this manner, a set of policers may be applied to the Voice VLAN to ensure that voice, video, and call signaling traffic are performing within normal levels and a separate, more stringent, policer can be applied to the data VLAN.



Note

When configuring policers for TelePresence, make sure you allow for the appropriate burst intervals, as defined in [Burst Requirements](#) in Chapter 4, "Quality of Service Design for TelePresence."

Finally, to ensure guaranteed levels of service, queuing needs to be configured on all nodes where the potential for congestion exists, regardless of how infrequently it may occur.

In Catalyst switches, queuing (along with all other QoS operations) is performed in hardware. Therefore, there are a fixed number of hardware queues that vary by platform, as well as by linecards. The nomenclature for Catalyst queuing is 1PxQyT, where:

- 1P represents a strict priority (Expedite Forwarding) queue
- xQ represents a number of non-priority queues
- yT represents a number of drop-thresholds per queue

**Note**

As discussed, due to the strict service levels required by Cisco TelePresence, it is recommended to assign TelePresence flows to a strict priority queue, whether this is implemented in Cisco Catalyst hardware or in Cisco IOS software. However, some older Catalyst platforms and linecards do not support a strict priority queue. For example, some Catalyst 6500 linecards support only a 2Q2T egress queuing model and as such would not be recommended within a Cisco TelePresence campus network design.

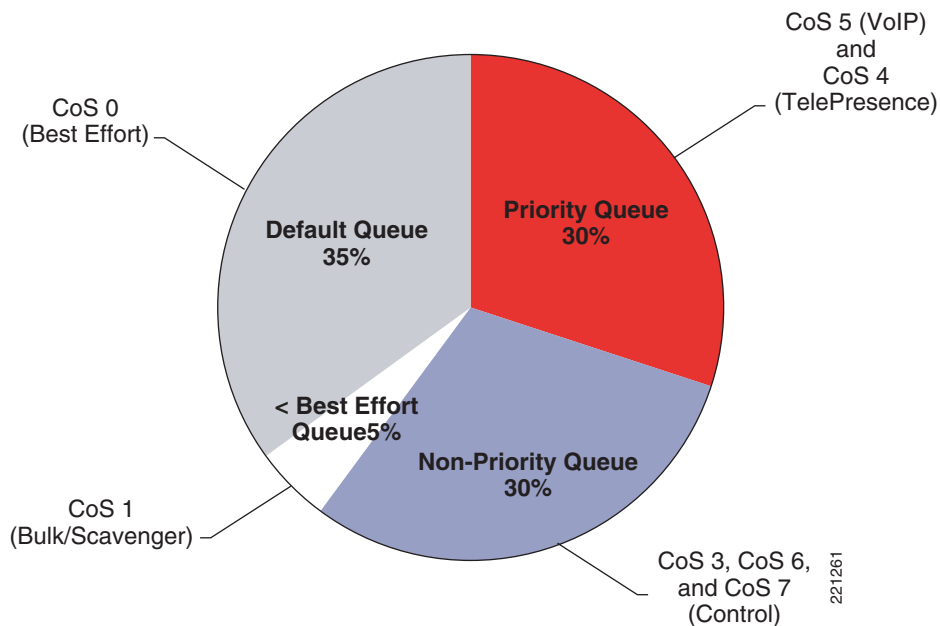
It is highly recommended that all Catalyst switches and linecards within a Cisco TelePresence campus design support a 1PxQyT queuing model.

For example, a Catalyst 6500 48-port 10/100/1000 RJ-45 Module (WS-X6748-GE-TX) has a 1P3Q8T, meaning 1 strict priority queue (which, incidentally, on this linecard is Queue 4) and 3 additional non-priority queues each with 8 configurable Weighted Random Early Detect (WRED) drop thresholds per queue.

Cisco Enterprise Systems Engineering (ESE) testing has shown that the optimal and most consistent service levels for TelePresence are achieved when TelePresence is provisioned with strict-priority hardware queuing (typically in conjunction with VoIP Telephony traffic), provided that the total bandwidth assigned to these realtime applications is less than 33% of the link—but this is virtually always the case on high-speed campus links in the range of 100 Mbps to 10 Gbps Ethernet. For example, consider a Catalyst 6513 provisioned with 11 x 48-port linecards, with each port configured to support G.711 VoIP (128 kbps max per port). Such a configuration would only require 67.584 Mbps or 6.8% of a GigE uplink. Even if a CTS-3000 system were connected to each of the 11 linecards, the total realtime bandwidth would be $[(11 \times 15 \text{ Mbps}) + (11 \times 48 \times 128 \text{ kbps})]$ 232.584 Mbps or 23.3% of a GigE uplink (which is still within the 33% LLQ Rule allowance).

As a generic campus queuing guide, it would be recommended to assign CoS values 4 (TelePresence) and 5 (VoIP) to the strict priority queue, CoS 3 (Call-Signaling) to a (non-default) non-priority queue, and CoS 0 (Best Effort) to the default queue. Finally, CoS 1 (Bulk and/or Scavenger traffic) should be assigned to a (minimally provisioned) less than Best Effort non-priority queue. These guidelines are illustrated in [Figure 5-2](#).

Figure 5-2 Generic Campus Queuing Provisioning and Mapping Guidelines



Campus Inter-Switch Link QoS Considerations

Once the trust boundary has been established and optimal access-edge policers have been enabled, then the DSCP values on all other inter-switch links and campus-to-WAN hand off links can be trusted. Therefore, it is recommended to trust DSCP (not CoS) on all inter-switch links, whether these are uplinks/downlinks to/from the distribution layer, uplinks/downlinks to/from the core layer, intra-core links, or links to WAN Aggregation routers.

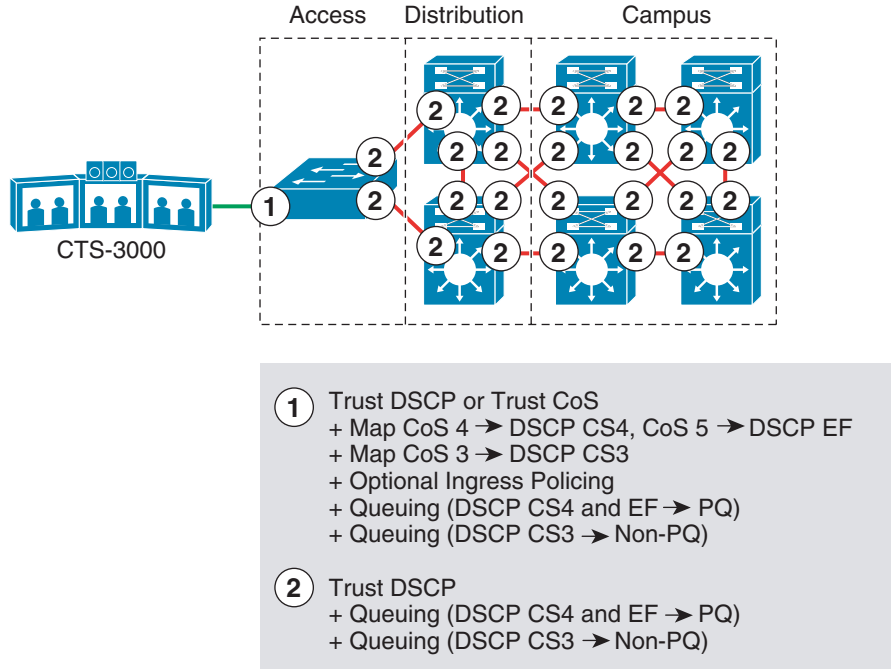
The reason it is recommended to trust DSCP and not CoS is two-fold: first, because marking granularity is lost every time a node is set to trust CoS. For example, if TelePresence endpoints are marking traffic to CS4 and Unified Video Advantage (or other Videoconferencing/Video Telephony endpoints) are marking their traffic to AF41 and the distribution-layer is set to trust CoS from the access-layer, then these flows both appear the same (as CoS 4) to the distribution-layer switch and are indistinguishable from each other from that node forward. Secondly, because trusting CoS implies using 802.1Q trunking between switches. Today, most enterprise campus networks are designed to be Layer 3 and thus 802.1Q is not used on inter-switch links.

Queuing is likewise recommended to be enabled on every node along the path. Note that this document generally focuses only on the QoS requirements for TelePresence. The actual QoS policies may be more complex than those shown here due to the myriad of other data, voice, and video applications on the network. It is recommended that customers use the information provided in this document in concert with

- *Enterprise QoS Solution Reference Network Design Guide*, Version 3.3, November 2005
- Szigeti, Tim and Hattingh, Christina. *End-to-End QoS Network Design: Quality of Service in LANs, WANs, and VPNs*. Indianapolis: Cisco Press, 2004. ISBN-10: 1-58705-176-1; ISBN-13: 978-1-58705-176-0.

A summary of the minimum QoS design requirements within an enterprise campus supporting TelePresence are illustrated in [Figure 5-3](#).

Figure 5-3 Enterprise Campus QoS Design Recommendations for TelePresence



TelePresence Campus QoS Designs

Now that campus-specific considerations have been addressed, we can look at how these policies can be configured on specific platforms. Before we identify the platforms, let's briefly review some of the key network and QoS-related features required by switch platforms at the access, distribution, and core layers of the enterprise campus. These include GE (including 10/100/1000) or 10GE connectivity, adequate port buffering to accommodate TelePresence traffic rates and micro-second bursts, granular policing, and IPxQyT queuing. Optionally, it would be preferred to have dedicated per-port buffers, as well as support for conditional trust, Per-Port/Per-VLAN policing, and DSCP-to-Queue mapping (as opposed to CoS-to-Queue mapping).

Given these requirements, the following currently-shipping Catalyst platforms have been validated by Cisco Enterprise Systems Technical Marketing for TelePresence-enabled campus networks include:

- Catalyst 3560G and 3750G (and by extension the 3650-E and 3750-E)
- Catalyst 4500 and 4948
- Catalyst 6500 (Although only certain linecards are recommended. Some older linecards do not have the requisite buffer to handle TelePresence traffic rate and burst requirements.)

Platform-specific configurations for each of these series of switches are provided in subsequent sections.

Catalyst 3560G/3750G and 3650-E/3750E

The Cisco Catalyst 3560G is a fixed-configuration switch that supports up to 48 10/100/1000 ports with integrated Power over Ethernet (PoE), plus 4 Small Form-Factor Pluggable (SFP) ports for uplinks. The 3560 has a 32 Gbps backplane, which is moderately oversubscribed (52 Gbps theoretical maximum input vs. 32 Gbps backplane yields an oversubscription ratio of 1.625:1 or 13:8). Additionally, the 3560G supports IP routing (including IPv6), multicast routing, and an advanced QoS and security feature-set.

The Catalyst 3750G is nearly identical, with only a few additional key features, including the support for a stackable configuration (via Stackwise technology), allowing for the 32 Gbps backplane (comprised of dual counter-rotating 16 Gbps rings) to be extended over multiple 3750G switches (up to 9). Additionally, the 3750G provides support for 10 Gigabit Ethernet (10GE) connectivity. Obviously, however, the more switches in the stack, as well as the use of 10 GE connectors, increases the oversubscription ratio accordingly.

The 3560-E and 3750-E represent the next evolution of these switches. As before, the 3560-E is a fixed configuration switch, but now with a 128 Gbps backplane and 10 GE port support. Similarly, the 3750-E supports a 128 Gbps backplane with dual 10GE port support, as well as the support for a stackable configuration (via Stackwise Plus technology, allowing a 64 Gbps interconnect between stacked switches).

As the 3560G, 3750G, 3560-E, and 3750-E share virtually identical feature parity (the main differences being the backplane throughput and uplink port speeds), we consider them as a single switch and abbreviate the reference to simply C3560G/3750G.

From a QoS perspective, some of the relevant features of the C3560G/3750G/E include conditional trust, Per-Port/Per-VLAN policers (via Hierarchical QoS policies), DSCP-to-Queue mapping, 2Q3T or 1P1Q3T ingress queuing, and 4Q3T or 1P3Q3T egress queuing. Additionally, these platforms provide (minimally) 750 KB of receive buffers and 2 MB of transmit buffers for each set of 4 ports. These buffers can be allocated, reserved, or dynamically borrowed from a common pool, on a port-port, per-queue basis, depending on the administrative configurations chosen.

Let's begin leveraging these features into the validated best-practice designs for this switch family for supporting TelePresence at the campus access-layer.

As QoS is disabled by default on these switches, the first step that we must take is to globally enable QoS. We can do this by issuing the global command:

```
mls qos
```

With QoS enabled, we can configure the access-edge trust boundaries. As discussed previously, we have three options: trust DSCP, trust CoS, or conditional trust. It is recommended that ports used for data and VoIP Telephony be configured to conditionally trust CoS, while ports used for TelePresence be configured to either trust DSCP, trust CoS or conditionally trust CoS. Trusting DSCP on these ports is the simplest operationally. The interface command to configure DSCP trust is fairly straightforward:

```
mls qos trust dscp
```



Note

While Cisco IOS allows the configuration of trust CoS and conditional trust on uplink ports, uplink ports should be set to trust DSCP (only). This is required on the C3560G/3750G/E for two reasons: first, to preserve marking granularity between switches (as previously discussed in [Campus Inter-Switch Link QoS Considerations](#)), as well as to activate the DSCP-to-Queue mapping (versus the CoS-to-Queue mapping) on the uplink switch ports.

If you choose to trust CoS or conditionally trust CoS, ensure that the fifth parameter in the global CoS-to-DSCP map—which corresponds to the DSCP mapping for CoS 4—is set to 32 (CS4). Additionally, to support IP Telephony properly, the global CoS-to-DSCP mapping table should be modified such that CoS 5 (the sixth parameter in the CoS-to-DSCP map) is mapped to 46 (EF)—which is not the default (the default setting is 40/CS5). These settings are achieved via the following global and interface commands:

```
mls qos map cos-dscp 0 8 16 24 32 46 48 56
interface Gigx/y
  mls qos trust cos
```

If you choose to implement conditional trust on the TelePresence ports, it can be enabled with the following interface command:

```
mls qos trust device cisco-phone
```


Note

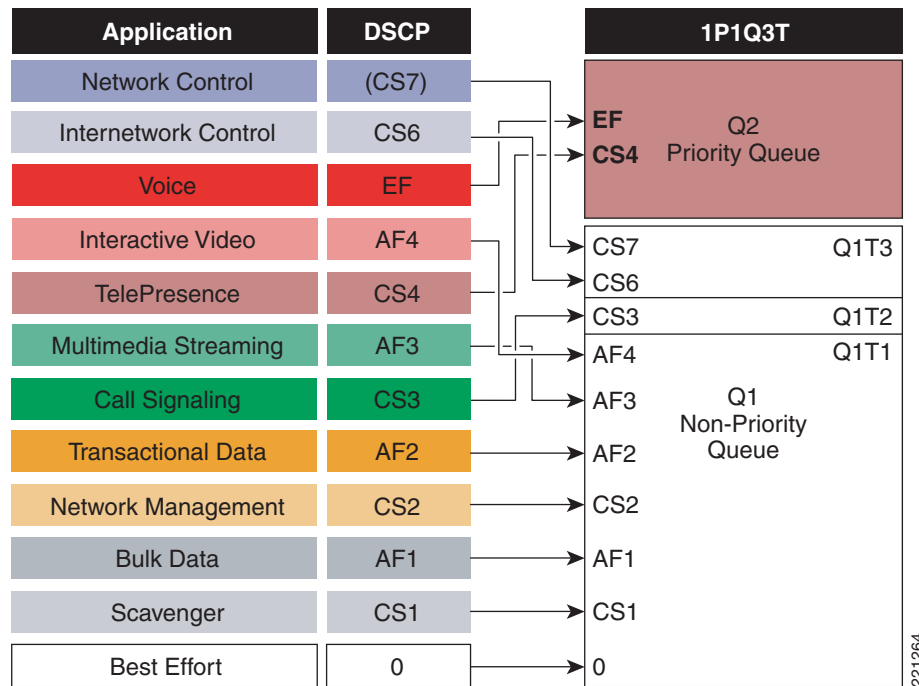
If conditional trust policies are to be used, then make sure that the TelePresence codec software is running version 1.1.0 (256D) or higher, as software version 1.0.1 (616D) incorrectly marks TelePresence audio traffic to CoS 5 (not CoS 4).

These configuration commands can be verified with the following commands:

- **show mls qos**
- **show mls qos map cos-dscp**
- **show mls qos interface**

Next, as the C3560G/3750G/E platforms have architectures based on oversubscription, they have been engineered to guarantee QoS by protecting critical traffic trying to access the backplane/stack-ring via ingress queuing. Ingress queuing on this platform can be configured as 2Q3T or 1P1Q3T. As we've already established the requirement for strict-priority servicing of TelePresence (and VoIP) traffic, it is recommended to enable the 1P1Q3T ingress queuing structure with DSCP EF (VoIP) and CS4 (TelePresence) being mapped to the ingress PQ (Q2). The configurable thresholds in the non-priority queue can be used to protect control traffic. For example, Network Control traffic (such as Spanning Tree Protocol) associated with DSCP CS7 and Internetwork Control traffic (such as Interior Gateway Protocols, including EIGRP and OSPF) marked DSCP CS6 can be explicitly protected by assigned these to Q1T3. Additionally, a degree of protection can be offered to Call-Signaling traffic (which is essentially control traffic for the IP Telephony infrastructure), which is marked CS3. All other traffic types can be provisioned in Q1T1. The recommended ingress 1P1Q3T queuing configuration for the C3560G/3750G/E platforms is illustrated in [Figure 5-4](#).

Figure 5-4 Catalyst 3560G/3750G/E(1P1Q3T) Ingress Queuing Recommendations for TelePresence Deployments



Based on [Figure 5-4](#), the recommended configuration for ingress queuing on the C3560G/3750G/E for TelePresence deployments is as follows:

! This first section modifies the CoS-to-DSCP for VoIP

```
mls qos map cos-dscp 0 8 16 24 32 46 48 56
! Modifies CoS-to-DSCP mapping to map CoS 5 to DSCP EF
```

! This section configures the Ingress Queues and Thresholds for 1P1Q3T

```
mls qos srr-queue input buffers 70 30
! Configures the Ingress Queue buffers such that Q2 (PQ) gets 30% of buffers
mls qos srr-queue input priority-queue 2 bandwidth 30
! Configures the Ingress PQ (Q2) to be guaranteed 30% BW on stack ring
mls qos srr-queue input bandwidth 70 30
! Configures SRR weights between Ingress Q1 and Q2 for remaining bandwidth
mls qos srr-queue input threshold 1 80 90
! Configures Ingress Queue 1 Threshold 1 to 80% and Threshold 2 to 90%
! Ingress Queue 1 Threshold 3 remains at 100% (default)
! Ingress Queue 2 Thresholds 1, 2 and 3 remain at 100% (default)
```

! This section configures the Ingress CoS-to-Queue Mappings for TelePresence ports using trust-CoS

```
mls qos srr-queue input cos-map queue 1 threshold 1 0 1 2
! Maps CoS 0, 1, 2 and 4 to Ingress Queue 1 (Q1T1)
mls qos srr-queue input cos-map queue 1 threshold 2 3
! Maps CoS 3 to Ingress Queue 1 Threshold 2 (Q1T2)
mls qos srr-queue input cos-map queue 1 threshold 3 6 7
! Maps CoS 6 and 7 to Ingress Queue 1 Threshold 3 (Q1T3)
mls qos srr-queue input cos-map queue 2 threshold 1 4 5
! Maps CoS 4 (TelePresence) and CoS 5 (VoIP) to Ingress-PQ Threshold 1 (Q2T1)
```

```

! This section configures the Ingress DSCP-to-Queue Mappings for TelePresence ports using
trust-DSCP

mls qos srr-queue input dscp-map queue 1 threshold 1 0 8 10 12 14
! Maps DSCP 0, CS1 and AF1 to Ingress Queue 1 Threshold 1 (Q1T1)
mls qos srr-queue input dscp-map queue 1 threshold 1 16 18 20 22
! Maps DSCP CS2 and AF2 to Ingress Queue 1 Threshold 1 (Q1T1)
mls qos srr-queue input dscp-map queue 1 threshold 1 26 28 30 34 36 38
! Maps DSCP AF3 and AF4 to Ingress Queue 1 Threshold 1 (Q1T1)
mls qos srr-queue input dscp-map queue 1 threshold 2 24
! Maps DSCP CS3 to Ingress Queue 1 Threshold 2 (Q1T2)
mls qos srr-queue input dscp-map queue 1 threshold 3 48 56
! Maps DSCP CS6 and CS7 to Ingress Queue 1 Threshold 3 (Q1T3)
mls qos srr-queue input dscp-map queue 2 threshold 1 32 46
! Maps DSCP CS4 (TelePresence)& EF (VoIP) to Ingress-PQ Threshold 1 (Q2T1)

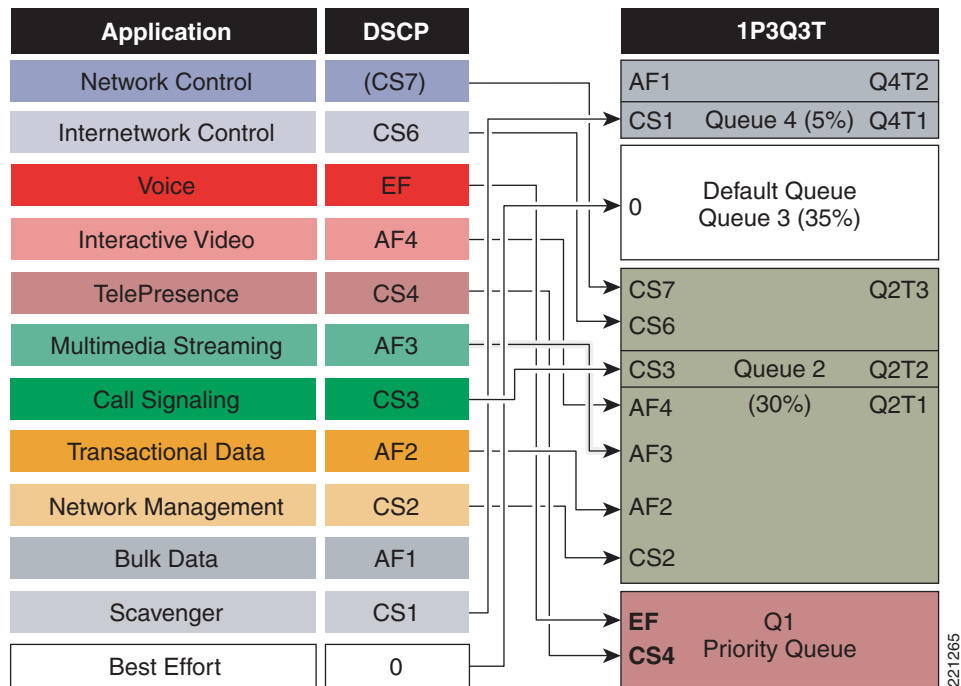
```

**Note**

Non-Standard DSCP values can also be mapped to their respective queues (using the CoS-to-Queue Map as a reference); however, for the sake of simplicity, non-standard DSCP-to-Queue Mappings have not been shown in these configurations.

Following ingress queuing configuration, we can now proceed to configuring the egress queues. The C3560G/3750G/E supports either 4Q3T or 1P3Q3T egress queuing configurations. As the need for an EF PHB has already been established, both for VoIP and for TelePresence, it is recommended to enable the 1P3Q3T egress queuing configuration, with Q1 as the PQ. Then both VoIP (DSCP EF) and TelePresence (DSCP CS4) should be mapped to Q1 (the PQ). Default traffic can be assigned to Q3 and Q4 can be designated as a less than Best Effort queue, servicing Bulk (AF1) and Scavenger (DSCP CS1) traffic, being assigned to Q4T2 and Q4T1, respectively. Network Control (DSCP CS7) and Internetwork Control (DSCP CS6) can be mapped to the highest threshold of the preferential non-priority queue (Q2T3), while Call-Signaling (DSCP CS3) can be mapped to the second highest threshold in that queue (Q2T2). All other applications can be mapped to Q2T1. The recommended 1P3Q3T egress queuing configuration for the C3560G/3750G/E platforms is illustrated in [Figure 5-5](#).

Figure 5-5 Catalyst C3560G/3750G/E (1P3Q3T) Egress Queuing Recommendations for TelePresence Deployments



Based on [Figure 5-5](#), the recommended configuration for egress queuing on the C3560G/3750G/E for TelePresence deployments is as follows:

! This section configures the Output CoS-to-Queue Maps for TelePresence ports using trust-CoS

```

mls qos srr-queue output cos-map queue 1 threshold 3 4 5
! Maps CoS 4 (TelePresence) and CoS 5 (VoIP) to Egress Queue 1 Threshold 3 (PQ)
mls qos srr-queue output cos-map queue 2 threshold 1 2
! Maps CoS 2 to Egress Queue 2 Threshold 1 (Q2T1)
mls qos srr-queue output cos-map queue 2 threshold 2 3
! Maps CoS 3 (Call-Signaling) to Egress Queue 2 Threshold 2 (Q3T2)
mls qos srr-queue output cos-map queue 2 threshold 3 6 7
! Maps CoS 6 and CoS 7 (Net Control) to Egress Queue 2 Threshold 3 (Q2T3)
mls qos srr-queue output cos-map queue 3 threshold 3 0
! Maps CoS 0 (Best Effort) to Egress Queue 3 Threshold 3 (Q3T3)
mls qos srr-queue output cos-map queue 4 threshold 3 1
! Maps CoS 1 (Bulk/Scavenger) to Egress Queue 4 Threshold 3 (Q4T3)

```

! This section configures the Output DSCP-to-Queue Maps for TelePresence ports using trust-DSCP

```

mls qos srr-queue output dscp-map queue 1 threshold 3 32 46
! Maps DSCP CS4 (TelePresence) and EF (VoIP) to Egress Queue 1 (PQ)
mls qos srr-queue output dscp-map queue 2 threshold 1 16 18 20 22
! Maps DSCP CS2 and AF2 to Egress Queue 2 Threshold 1 (Q2T1)
mls qos srr-queue output dscp-map queue 2 threshold 1 26 28 30 34 36 38
! Maps DSCP AF3 and AF4 to Egress Queue 2 Threshold 1 (Q2T1)
mls qos srr-queue output dscp-map queue 2 threshold 2 24
! Maps DSCP CS3 to Egress Queue 2 Threshold 2 (Q2T2)
mls qos srr-queue output dscp-map queue 2 threshold 3 48 56
! Maps DSCP CS6 and CS7 to Egress Queue 2 Threshold 3 (Q2T3)
mls qos srr-queue output dscp-map queue 3 threshold 3 0

```

```

! Maps DSCP DF to Egress Queue 3 Threshold 3 (Q3T3 - Default Queue)
mls qos srr-queue output dscp-map queue 4 threshold 1 8
! Maps DSCP CS1 to Egress Queue 4 Threshold 1 (Q4T1)
mls qos srr-queue output dscp-map queue 4 threshold 2 10 12 14
! Maps DSCP AF1 to Egress Queue 4 Threshold 2 (Q4T2)

! This next section configures the WRED min and max thresholds for Q1

mls qos queue-set output 1 threshold 2 80 90 100 100
! Sets Egress Queue 2 Threshold 1 (Q2T1) to 80% and Threshold2 (Q2T2) to 90%
mls qos queue-set output 1 threshold 4 60 100 100 100
! Sets Egress Queue 4 Threshold 1 (Q4T1) to 60% and Threshold 2 (Q4T2) to 100%

! This section configures trust-DSCP and queuing on TelePresence access port and uplink
ports

interface GigabitEthernet1/0/1
description TelePresence or Uplink port
mls qos trust dscp
! Assigns the TelePresence port and/or uplink port to trust DSCP
queue-set 1
! Assigns interface to Queue-Set 1 (default)
srr-queue bandwidth share 1 30 35 5
! Q2 gets 30% of remaining BW (after PQ); Q3 gets 35% & Q4 gets 5%
priority-queue out
! Expedite queue is enabled for TelePresence and VoIP
!

! This section configures conditional-trust and queuing on TelePresence access ports

interface GigabitEthernet1/0/2
description IP Telephony and/or Data port
mls qos trust device cisco-phone
! Configures conditional trust based on the CDP advertisements of the TelePresence
system and attached 7975G IP phone
queue-set 1
! Assigns interface to Queue-Set 1 (default)
srr-queue bandwidth share 1 30 35 5
! Q2 gets 30% of remaining BW (after PQ); Q3 gets 35% & Q4 gets 5%
priority-queue out
! Expedite queue is enabled for TelePresence and VoIP
!

```

**Note**

As before, non-Standard DSCP values can also be mapped to their respective queues (using the CoS-to-Queue Map as a reference); however, for the sake of simplicity, non-standard DSCP-to-Queue Mappings have not been shown in these configurations.

These configuration commands can be verified with the following commands:

- **show mls qos queue-set**
- **show mls qos maps cos-input-q**
- **show mls qos maps dscp-input-q**
- **show mls qos maps cos-output-q**
- **show mls qos maps dscp-output-q**

- **show mls qos interface**
- **show mls qos interface buffers**
- **show mls qos interface queueing**
- **show controllers ethernet-controller port-asic statistics**

Catalyst 4500 and 4948

The Cisco Catalyst 4500 series switches are midrange modular platforms with chassis options to support 3, 6, 7, and 10 slots; these models include the Catalyst 4503, 4506, 4507R, and 4510R, respectively (the latter two models supporting a redundant supervisor option). The Catalyst 4500 family of switches provides Layer 2 through Layer 4 network services, including advanced high-availability, security, and QoS services in addition to integrated PoE to support unified communications. The linecards that meet the requirements (at the time of writing) outlined in [TelePresence Campus QoS Designs](#) for the Catalyst 4500 include the 4448 and the 4548 series linecards (specifically, the WS-X4448-GB-RJ45 and the WS-X4524-GB-RJ45V or WS-X4548-GB-RJ45V).

On the other hand, the degree of oversubscription and buffering capabilities on the C4500 series linecards varies by linecard. Some linecards are entirely non-blocking, while others, such as the 4448 and the 4548, provision a single 1 Gbps uplink to the switch fabric for every 4 or 8 (10/100/1000) ports, which equates to an 4:1 (for the 4524) or an 8:1 (4448 and 4548) theoretical oversubscription ratio. As such, the 4448 and 4548 series linecards, while suitable at the campus access-edge, would not be recommended to be used as uplinks nor within the distribution and core layers of a TelePresence-enabled campus.

The Catalyst 4948 series provides an advanced feature set of intelligent network services, but is engineered and optimized to support high-performance wire-speed switching for data center server traffic. Thus the Catalyst 4948 has a completely non-blocking architecture and as such would be suitable at any layer (access, distribution, or core) within a TelePresence-enabled campus network. Specifically, the Catalyst 4948 provides 96 Gbps of switching fabric for its fixed configuration 48 x 10/100/1000 ports plus 4 SFP ports (which may be GE or 10 GE). Additionally, the Catalyst 4948 provides approximately 16 MB of buffering which is shared among all 48 ports.

As the Catalyst 4500 and 4948 share virtually identical feature parity (the main differences being the backplane throughput and buffer architectures), we consider them as a single switch and abbreviate the reference to simply C4500/4948.

From a QoS perspective, some of the relevant features of the C4500/4948 include conditional trust, an elegant Per-Port/Per-VLAN policer implementation, DSCP-to-Queue mapping, 4Q1T or 1P3Q1T queuing support, and an advanced congestion algorithm (Dynamic Buffer Limiting or DBL).

Let's begin leveraging these features into the validated best-practice designs for this switch family for supporting TelePresence at the campus access-layer.

The first thing to note is a minor syntactical difference when configuring QoS features on the C4500/4948; specifically, QoS commands on this platform do not include the mls prefix used on the C3560G/3750G/E and the C6500 series platforms. For example, to globally enable QoS on the C4500/4948 (which is disabled by default), the command is not mls qos, but simply:

```
qos
```

With QoS enabled, we can configure the access-edge trust boundaries. As discussed previously, we have three options: trust DSCP, trust CoS, or conditional trust. It is recommended that ports used for data and VoIP Telephony be configured to conditionally trust CoS, while ports used for TelePresence be configured to either trust DSCP, trust CoS or conditionally trust CoS. Trusting DSCP on these ports is the simplest operationally.

```
qos trust dscp
```

If you choose to trust CoS or conditionally trust CoS, then CoS 5 must be explicitly mapped to DSCP EF prior to the port being configured to trust CoS. All other CoS-to-DSCP mappings can be left at their respective default values. These functions can be achieved via the following global and interface commands:

```
qos map cos 5 to 46
!
interface Gigx/y
  qos trust cos
```

If you choose to implement conditional trust on the TelePresence ports, it can be enabled with the following interface command:

```
qos trust device cisco-phone
```


Note

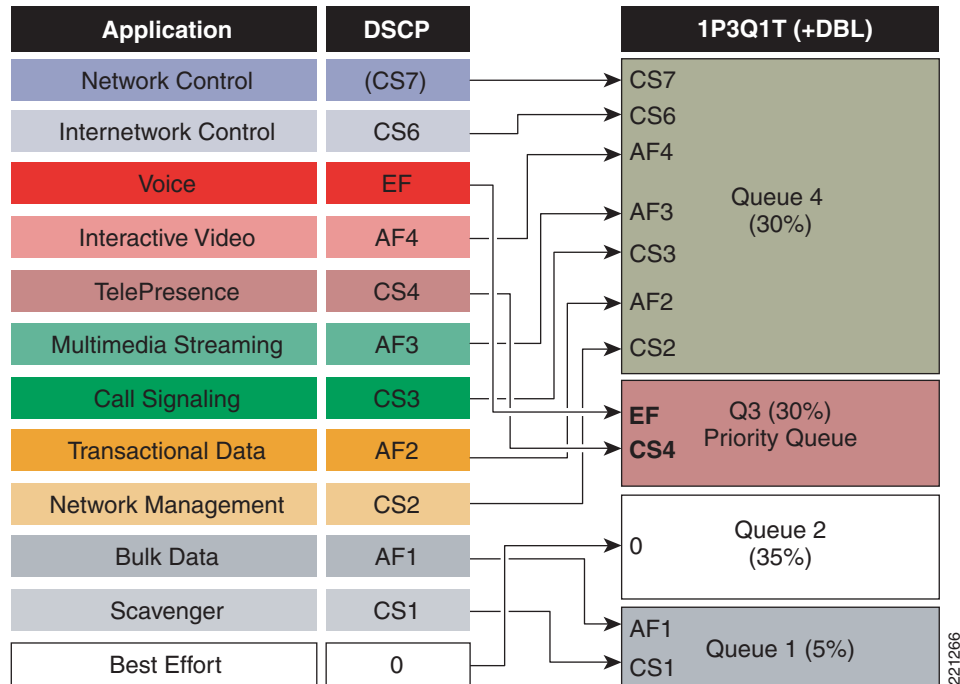
If conditional trust policies are to be used, then make sure that the TelePresence codec software is running version 1.1.0 (256D) or higher, as software version 1.0.1 (616D) incorrectly marks TelePresence audio traffic to CoS 5 (not CoS 4).

As with configuration commands, the C4500/4948 omits the mls prefix in the corresponding verification commands. These configuration commands can be verified with the following commands:

- **show qos**
- **show qos maps**
- **show qos interface**

As the C4500/4948 does not support ingress queuing (although it bears mentioning that the internal servicing architectures have been tested and found to be adequate in protecting TelePresence traffic even in the event of oversubscription), we can move on to configuring egress queuing. The C4500/4948 can be configured to operate in a 4Q1T mode or a 1P3Q1T mode, the latter of which is recommended for VoIP Telephony and TelePresence deployments. On the C4500, however, the strict priority queue, when enabled, is Q3. As the C4500/4948 supports DSCP-to-Queue mappings, we can distinguish between applications such as generic Videoconferencing/Video Telephony (AF4) and TelePresence (CS4), even though these share the same CoS and IP Precedence values (Cos/IPP 4). Given these abilities, it is recommended to enable 1P3Q1T queuing on the C4500/4948, with VoIP (EF) and TelePresence (CS4) assigned to the strict-priority queue (Q3). Q2 may be dedicated to service default traffic and Q1 can be used to service less than Best Effort Scavenger (CS1) and Bulk (AF1) traffic. All other applications can be mapped to Q4, the preferential queue. The recommended (1P3Q1T + DBL) egress queuing configuration for the C4500/4948 platform is illustrated in [Figure 5-6](#).

Figure 5-6 Catalyst C4500/4948 (1P3Q1T + DBL) Egress Queuing Recommendations for TelePresence Deployments



Note

As before, non-Standard DSCP values can also be mapped to their respective queues; however, for the sake of simplicity, non-standard DSCP-to-Queue Mappings have not been shown in these configurations.

As previously mentioned, the C4500/4948 supports an advanced congestion avoidance algorithm—Dynamic Buffer Limiting (DBL)—rather than Weighted Tail Drop (WTD) or Weighted-Random Early-Detect (WRED). Therefore no DSCP-to-Threshold mappings are required on the C4500/4948. However, to leverage DBL, it must be globally enabled (as it is disabled by default). This is achieved with the following global command:

```
qos db1
```

Optionally, DBL can be configured to operate to support RFC 3168 IP Explicit Congestion Notification (IP ECN or simply ECN), which utilizes the remaining 2 bits of the IPv4/IPv6 Type of Service (ToS) Byte (the DSCP value uses the first 6 bits of the ToS Byte). The following global command enables ECN for DBL:

```
qos db1 exceed-action ecn
```



Note

For more information on IP ECN, refer to RFC 3168 (at www.ietf.org/rfc/rfc3168) and Sziget, Tim and Hattingh, Christina. *End-to-End QoS Network Design: Quality of Service in LANs, WANs, and VPNs*. Indianapolis: Cisco Press, 2004. ISBN-10: 1-58705-176-1; ISBN-13: 978-1-58705-176-0.

Additionally, to leverage DBL (with/without ECN) on a per-interface basis, a service policy applying DBL to all flows must be constructed and applied to each interface. This can be done by using the following basic policy-map:

```
policy-map DBL
```

```

class class-default
  db1
!
interface Gig x/y
  service policy output DBL

```

However, at this point, an important consideration pertaining to DBL must be taken into account, namely DBL (when enabled and configured as per the above recommendations) is active on all flows, including flows destined to the PQ (Q3)—which in our case includes VoIP and TelePresence traffic. As DBL introduces dynamic drops, especially on bursty, large-packet flows, this is detrimental to TelePresence call-quality. Therefore, to explicitly disable DBL on PQ traffic, the following amendments can be made to the previous policy:

```

class-map PQ
  match ip dscp ef
  match ip dscp cs4
policy-map DBL
  class PQ
  class class-default
    db1
!
interface Gig x/y
  service policy output DBL

```

In this modified policy, the class-map PQ identifies traffic destined to the Priority Queue, specifically EF (VoIP) and CS4 (TelePresence) traffic. In the policy-map, the PQ-class receives no action (DBL or otherwise) and serves only to exclude these flows from the following class-default policy of applying DBL to all (other) flows. It is highly recommended to use this modified policy on C4500/4948 platforms supporting TelePresence in conjunction with DBL; otherwise DBL drops negatively impact TelePresence call-quality.

Piecing this together, the C4500/4948 egress queuing recommendation, shown in [Figure 5-6](#), is as follows:

```

!This section enables DBL globally and excludes DBL on PQ flows

qos db1
  ! Globally enables DBL
qos db1 exceed-action ecn
  ! Optional: Enables DBL to mark RFC 3168 ECN bits in the IP ToS Byte
class-map PQ
  match ip dscp ef
  match ip dscp cs4
  ! Classifies traffic mapped to PQ for exclusion of DBL-policy
policy-map DBL
  class PQ
    ! No action (DBL or otherwise) is applied on traffic mapped to PQ
  class class-default
    db1
    ! Enables DBL on all (other) traffic flows

! This section configures the DSCP-to-Transmit Queue Mappings

qos map dscp 0 to tx-queue 2
  ! Maps DSCP 0 (Best Effort) to Q2
qos map dscp 8 10 12 14 to tx-queue 1
  ! Maps DSCP CS1 (Scavenger) and AF11/AF12/AF13 (Bulk) to Q1
qos map dscp 16 18 20 22 to tx-queue 4
  ! Maps DSCP CS2 (Net-Mgmt) and AF21/AF22/AF23 (Transactional) to Q4
qos map dscp 24 26 28 30 to tx-queue 4
  ! Maps DSCP CS3 (Call-Sig) and AF31/AF32/AF33 (MultiMedia) to Q4

```

```

qos map dscp 34 36 38 to tx-queue 4
! Maps DSCP AF41/AF42/AF43 (Interactive-Video) to Q4
qos map dscp 32 46 to tx-queue 3
! Maps DSCP CS4 (TelePresence) and EF (VoIP) to Q3 (PQ)
qos map dscp 48 56 to tx-queue 4
! Maps DSCP CS6 (Internetwork) and CS7 (Network Control) to Q4

! This section configures queues, activates the PQ and applies DBL

interface range GigabitEthernet1/1 - 48
  tx-queue 1
  bandwidth percent 5
  ! Q1 gets 5% BW
  tx-queue 2
  bandwidth percent 35
  ! Q2 gets 35% BW
  tx-queue 3
  priority high
  ! Q3 is PQ
  bandwidth percent 30
  ! Q3 (PQ) gets 30% BW
  shape percent 30
  ! Shapes/limits PQ to 30% BW
  tx-queue 4
  bandwidth percent 30
  ! Q4 gets 40%
  service-policy output DBL
  ! Applies DBL to all flows except VoIP & TelePresence
!
```

**Note**

As before, non-Standard DSCP values can also be mapped to their respective queues; however, for the sake of simplicity, non-standard DSCP-to-Queue Mappings have not been shown in these configurations.

These configuration commands can be verified with the following commands:

- **show qos dbl**
- **show qos maps dscp tx-queue**
- **show qos interface**

Catalyst 6500

The Cisco Catalyst 6500 series switches represent the flagship of Cisco's switching portfolio, delivering innovative secure, converged services throughout the campus, from the access-edge wiring closet to the distribution to the core to the data center to the WAN edge. The Catalyst 6500 platform is available in 3, 4, 6, 9, or 13 slot combinations; these models include the 6503, 6504, 6506, 6509 (regular or Network Equipment Building System [NEBS] compliant), and 6513. Additionally, these chassis options are also available in Enhanced models, designated by a -E suffix (such as 6503-E, 6504-E, etc.) for additional feature functionality and performance (except the 6513 at the time of writing).

Overall, the Catalyst 6500 provides the highest performance switching plane, supporting a 720 Gbps switching fabric and the option to run either centralized or distributed forwarding to achieve optimal performance. Additionally, the Catalyst 6500 provides leading-edge Layer 2-Layer 7 services, including rich High-Availability, Manageability, Virtualization, Security, and QoS feature sets, as well as integrated Power-over-Ethernet (PoE), allowing for maximum flexibility in virtually any role within the campus.

The queuing and buffering capabilities of Catalyst 6500 supervisors and linecards vary according to type. A summary of the ingress and egress queuing structures, as well as ingress, egress, and total buffering capabilities for Catalyst 6500 supervisors and linecards that have been tested and validated for TelePresence campus network designs are presented in [Figure 5-7](#) through [Figure 5-9](#).

[Figure 5-7](#) presents queuing and buffering details by Catalyst 6500 supervisor types.

Figure 5-7 Catalyst 6500 Queuing and Buffering Details by Supervisor Types

Supervisor Engines	Ingress Queue and Drop Thresholds	Ingress Queue Scheduler	Egress Queue and Drop Thresholds	Egress Queue Scheduler	Total Buffer Size	Ingress Buffer Size	Egress Buffer Size
WS-SUP720	1p1q4t	WRR	1p2q2t	WRR	512 KB	73 KB	439 KB
WS-SUP720-3B							
WS-SUP720-3BXL							
WS-SUP32-10GE	2q8t	WRR	1p3q8t	SRR	1.3 MB	166 KB	1.2 MB
WS-SUP32-GE							

224698

[Figure 5-8](#) presents queuing and buffering details for Catalyst 6500 GigabitEthernet linecards (including 10/100/1000 linecards).

Figure 5-8 Catalyst 6500 Queuing and Buffering Details by GigabitEthernet (and 10/100/1000) Linecards

Modules	Ingress Queue and Drop Thresholds	Ingress Queue Scheduler	Egress Queue and Drop Thresholds	Egress Queue Scheduler	Total Buffer Size	Ingress Buffer Size	Egress Buffer Size
WS-X6148A-GE-TX	1q2t	WRR	1p3q8t	WRR	5.5 MB	120 KB	5.4 MB
WS-X6148A-GE-45AF							
WS-X6516A-GBIC	1p1q4t	WRR	1p2q2t	WRR	1 MB	135 KB	946 /kB
WS-X6548-GE-TX	1q2t	WRR	1p2q2t	WRR	1.4 MB	185 KB	1.2 MB
WS-X6548V-GE-TX							
WS-X6548-GE-45AF							
WS-X6724-SFP with CFC	1q8t	WRR	1p3q8t	DWRR	1.3 MB	166 KB	1.2 MB
WS-X6724-SFP with DFC3	2q8t	WRR					
WS-X6748-GE-TX with CFC	1q8t	WRR					
WS-X6748-GE-TX with DFC3	2q8t	WRR					
WS-X6748-SFP with CFC	1q8t	WRR					
WS-X6748-SFP with DFC3	2q8t	WRR					

224699

It is important to note that the 6148A-GE and the 6548-GE are both engineered with 8:1 oversubscription ratios and, as such, while suitable at the access layer, these linecards would not be recommended to deploy as uplinks or within the distribution and core layers of the TelePresence-enabled campus network. Furthermore, it would be recommended to configure ingress queuing on these linecards, according to the recommendations presented in [Ingress Queuing Design—1Q2T](#).

In contrast, the 6748-GE is virtually non-blocking, supporting a dual 20 Gbps connection to the switch fabric for its 48 (10/100/1000) ports, which equates to a minimal 6:5 oversubscription ratio.

In turn, [Figure 5-9](#) presents queuing and buffering details for Catalyst 6500 TenGigabitEthernet (10GE) linecards.

Figure 5-9 Catalyst 6500 Queuing and Buffering Details by TenGigabitEthernet Linecards

Modules	Ingress Queue and Drop Thresholds	Ingress Queue Scheduler	Egress Queue and Drop Thresholds	Egress Queue Scheduler	Total Buffer Size	Ingress Buffer Size	Egress Buffer Size
WS-X6704-10GE with CFC	1q8t	WRR	1p7q8t	DWRR	16 MB	2 MB	14 MB
WS-X6704-10GE with DFC3	8q8t	WRR					
WS-X6708-10G-3C	8q4t	DWRR	1p7q4t	DWRR SRR	198 MB	108 MB	90 MB
WS-X6708-10G-3CXL							

224700

It is important to note that, at the time of writing, all Catalyst 6500 linecards support only CoS-to-Queue configurations for both ingress and egress queuing (with the exception of the 6708 [WS-X6708-10GE] linecard family, which can be configured with either CoS-to-Queue mappings or with DSCP-to-Queue mappings). Considerations relating to CoS-to-Queue mapping limitations are covered shortly.

Finally, [Figure 5-10](#) summarizes egress queuing structures by supervisor and linecard families, as a quick reference to selecting the relevant recommended queuing configurations.

Figure 5-10 Catalyst 6500 Supervisors and Linecards By Egress Queuing Structures

1P2Q2T CoS-to-Queue	1P3Q8T CoS-to-Queue	1P7Q8T CoS-to-Queue	1P7Q4T* *DSCP-to-Queue
<ul style="list-style-type: none"> • WS-SUP720 • WS-SUP720-3B • WS-SUP720-3BXL • WS-6516A-GBIC • WS-X6548-GE-TX • WS-X6548V-GE-TX • WS-X6548-GE-45AF 	<ul style="list-style-type: none"> • WS-SUP32-GE • WS-SUP32-10GE • WS-X6148A-GE-45F • WS-X6148A-GE-TX • WS-X6724-SFP • WS-X6748-GE-TX • WS-X6748-SFP 	<ul style="list-style-type: none"> • WS-X6704-10GE 	<ul style="list-style-type: none"> • WS-X6708-10G-3C • WS-X6708-10G-3CXL

224701

From a QoS perspective, some of the relevant features of the C6500 include port-trust, linecard-dependant queuing options, and WRED support. Let's examine how these features can be leveraged into validated best-practice designs for the Catalyst 6500 (which we will abbreviate to C6500) at the access-edge.

As with the previously discussed switch platforms, QoS is disabled by default and must be explicitly enabled globally on the C6500 for any configured policies to take effect. The command to globally enable QoS on the C6500 is:

```
mls qos
```

With QoS enabled, we can configure the access-edge trust boundaries. At the time of writing, on the C6500, we have only two port-trust options: trust DSCP and trust CoS.

**Note**

While trusting IP Precedence is a configurable option, this functionality is superseded by trusting DSCP. Additionally, at the time of writing, conditional trust is not available on the C6500.

When considering which trust option to configure, there is an important relationship between trust and ingress queuing on the C6500 to consider, namely, if a port is set to trust CoS, then ingress queuing is automatically enabled. This becomes an especially relevant consideration on linecards with high oversubscription ratios, such as the 6148A and 6548 (both with 8:1 oversubscription ratios). Therefore, it is recommended to set the ports connecting to TelePresence systems to trust CoS. However, keep in mind that if CoS is to be trusted, then ensure that the fifth parameter in the global CoS-to-DSCP map—which corresponds to the DSCP mapping for CoS 4—is set to 32 (CS4). Additionally, to support IP Telephony properly, the global CoS-to-DSCP mapping table should be modified such that CoS 5 (the sixth parameter in the CoS-to-DSCP map) is mapped to 46 (EF), which is not the default (the default setting is 40/CS5). These settings are achieved via the following global and interface commands:

```
mls qos map cos-dscp 0 8 16 24 32 46 48 56
interface Gigx/y
  mls qos trust cos
```

However, on all inter-switch link ports (uplinks/downlinks, etc.) it is recommended to set port-trust to trust DSCP to preserve marking granularity and Diffserv Per-Hop Behaviors. For this same reason, it is highly recommended to set port-trust to trust DSCP (or better yet, to use service policies with policers) on ports connected to endpoints that may be generating generic Videoconferencing/Video Telephony traffic (marked AF41); otherwise the DSCP value (AF41) for these flows will be lost and will be remapped to CS4 (since the CoS-to-DSCP map maps CoS 4 to DSCP CS4), eliminating your ability to distinguish between them on subsequent switch/router hops along the network path. The interface command to configure DSCP trust on an interface is as follows:

```
mls qos trust dscp
```

These configuration commands can be verified with the following commands:

- **show mls qos**
- **show mls qos maps | begin Cos-dscp map**
- **show queuing interface gigabitethernet x/y | include trust**

Before we describe linecard-specific queue recommendations, it bears mentioning that the queuing structures on the C6500—both ingress and egress—are CoS-based (with the sole exception, at the time of writing, of the WS-X6708-10GE, which supports DSCP-based queuing). This presents a challenge to network administrators deploying both generic Videoconferencing/Video Telephony (marked AF41 per RFC 4594) and TelePresence (marked CS4 per RFC 4594), as these applications both share the same CoS value of 4. As such, TelePresence and generic Videoconferencing/Video Telephony traffic are indistinguishable from one another with a CoS-based queuing scheme, with both applications always being mapped to the same queue. Since TelePresence requires an Expedited Forwarding Per-Hop Behavior (as explained in [Chapter 4, “Quality of Service Design for TelePresence”](#) and as allowed for by RFC 4594), both TelePresence and Videoconferencing/Video Telephony can be assigned to the strict-priority hardware queue on C6500 linecards (along with VoIP). Therefore, while this technical limitation exists, network administrators are encouraged to configure the hardware strict-priority queues on their C6500 platforms to adequately provision for their VoIP, TelePresence, and Videoconferencing/Video Telephony traffic.

While this may sound a bit complicated or excessive, in practice it is not that difficult to do, especially when considering that VoIP is such a lightweight application. For example, consider a Catalyst 6513 with redundant supervisors and 11 x 48-port linecards. If each of these ports supported VoIP, a total of only 68 Mbps of PQ would be required on the uplink (6.8% of a GE link). Additionally, if a generic

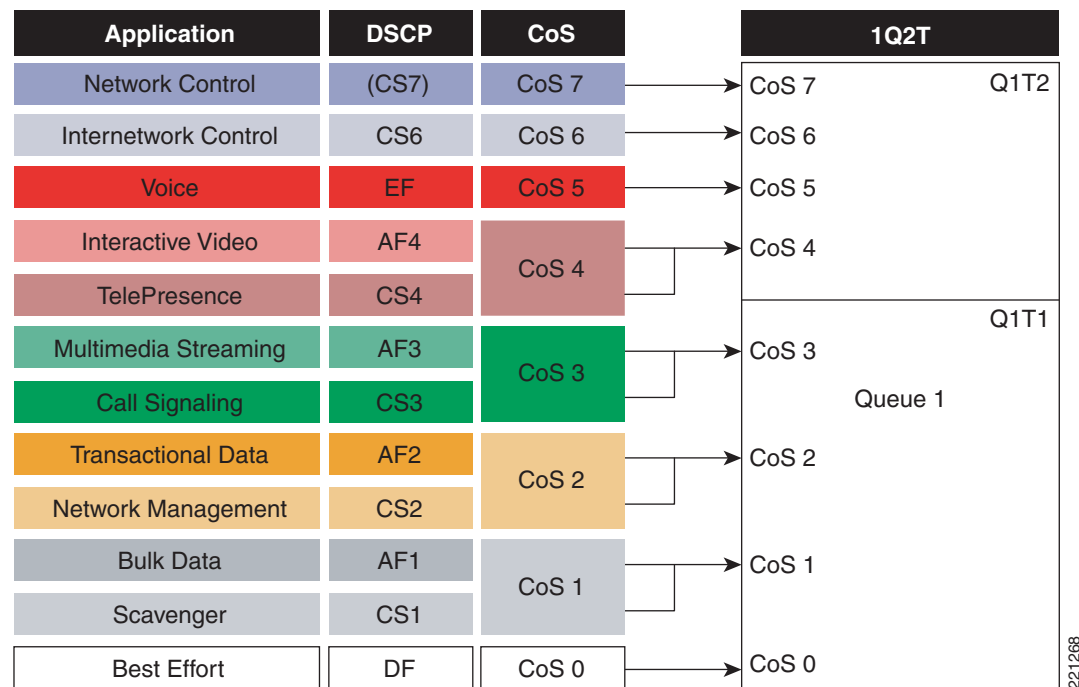
Videoconferencing/Video Telephony application was provisioned on each port that would permit 384 Kbps of AF41 video traffic per port, the combined total would be 270 Mbps (27% of a GE uplink). This leaves enough PQ traffic to support 2 separate CTS-3000 systems connected to the same chassis and still be at a theoretical maximum of only 30% of a single GE uplink.

With this in mind, let's now consider the best practice ingress and egress queuing configurations for the 6148A, 6548 and 6748 linecards.

Ingress Queuing Design—1Q2T

As shown in [Figure 5-10](#), both the 6148A and 6548 linecards support a CoS-based ingress queuing structure of 1Q2T which can be leveraged to offset their oversubscription ratios. The 1Q2T ingress queuing structure uses Tail-Drop thresholds, which by default are set at 80% of the queue (Q1T1) and at 100% of the queue (Q1T2 which, incidentally, is non-configurable). By default, CoS values 0 through 4 are mapped to Q1T1 and CoS values 5 through 7 are mapped to Q1T2. The only improvement we can make on this default configuration to optimize TelePresence traffic on these oversubscribed linecards is to map CoS 4 (TelePresence) to the second threshold (Q1T2), along with CoS 5 (VoIP) and CoS 6 and 7 (Network Control traffic). The recommending ingress 1Q2T queuing configuration for C6500 6148A and 6548 linecards is illustrated in [Figure 5-11](#).

Figure 5-11 Catalyst 6500 (1Q2T) Ingress Queuing Recommendations for TelePresence Deployments



The configuration for the C6500 1Q2T ingress queuing structure (for the 6148A and 6548 linecards) illustrated in [Figure 5-11](#) is as follows:

```
interface GigabitEthernet x/y
  rcv-queue cos-map 1 1 0 1 2 3
    ! Maps CoS values 0-3 to Q1T1
  rcv-queue cos-map 1 2 4 5 6 7
    ! Maps CoS values 4-7 to Q1T2
```

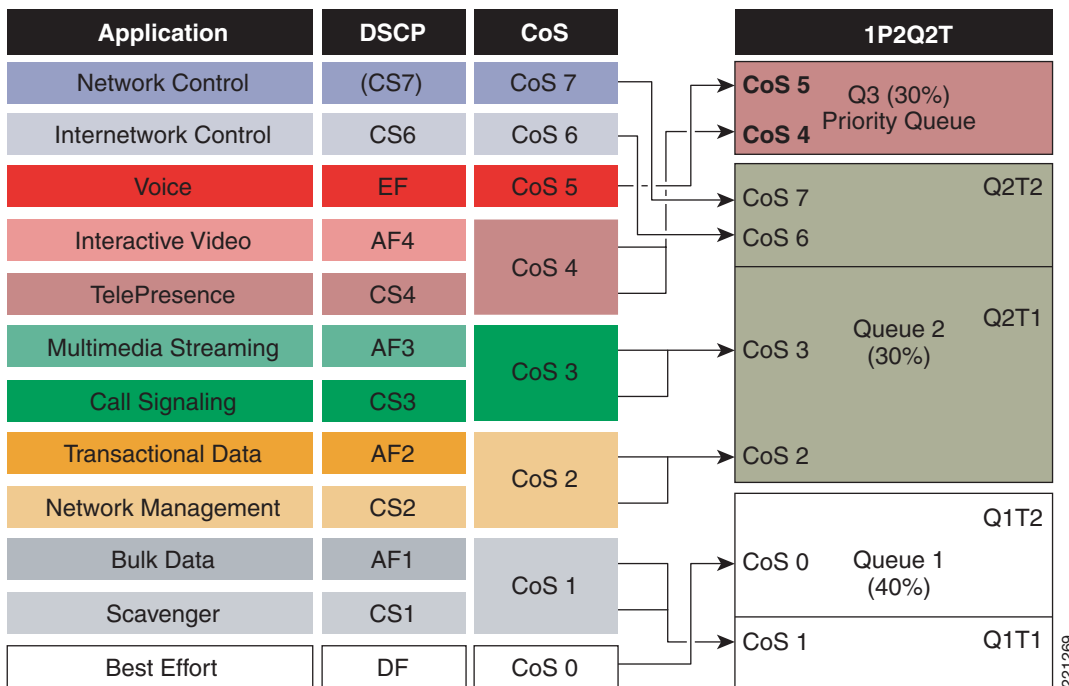
These configuration commands can be verified with the following command:

- **show queuing interface GigabitEthernet x/y**

Egress Queuing Design—1P2Q2T

As shown in [Figure 5-10](#), the Sup720s and 6548 linecards support a CoS-based egress queuing structure of 1P2Q2T, which uses WRED as a congestion avoidance mechanism. Under such a queuing structure, TelePresence traffic, along with VoIP, is recommended to be mapped to the strict-priority queue. Furthermore, because the egress queuing is CoS-based, Videoconferencing/Video Telephony (AF41) will also be assigned to the strict-priority queue. For non-realtime classes, the per-queue WRED thresholds can be configured to allow for granular QoS within a given queue. Specifically, Q1T1's minimum WRED threshold is set to 40% and its maximum threshold to 80%; then by mapping CoS 1 (Scavenger/Bulk) to Q1T1, we are restricting such traffic within Q1, with the remaining buffers (Q1T2) being exclusively reserved for CoS 0 (Best Effort traffic). Similarly, we can set Q2T1's minimum WRED threshold to 70% and maximum threshold to 80%; then by mapping CoS 2 (Transactional/Network Management) and CoS 3 (Call-Signaling/Multi-Media Streaming) to Q2T1, we are restricting these flows to a maximum of 80% of Q2, with the remaining buffers (Q2T2) being exclusively reserved for CoS 6 and 7 (Network Control traffic). The recommending egress 1P2Q2T queuing configuration for C6500 6548 linecards is illustrated in [Figure 5-12](#).

Figure 5-12 Catalyst 6500 (1P2Q2T) Egress Queuing Recommendations for TelePresence Deployments



The configuration for the C6500 1P2Q2T egress queuing structure (for the 6548 linecards) illustrated in [Figure 5-12](#) is as follows:

```
!
interface GigabitEthernet x/y

! This section sets the queue limits and bandwidth allocations
wrr-queue queue-limit 40 30
```

```

! Sets the buffer allocations to 40% for Q1 and 30% for Q2
! Also implicitly sets PQ (Q3) to 30%)
wrr-queue bandwidth 40 30
! Sets the WRR weights for 40:30 (Q1:Q2) bandwidth servicing

! This section sets the Min and Max WRED thresholds for Q1
wrr-queue random-detect min-threshold 1 40 80
! Sets Min WRED Thresholds for Q1T1 and Q1T2 to 40 and 80, respectively
wrr-queue random-detect max-threshold 1 80 100
! Sets Max WRED Thresholds for Q1T1 and Q1T2 to 80 and 100, respectively

! This section sets the Min and Max WRED thresholds for Q2
wrr-queue random-detect min-threshold 2 70 80
! Sets Min WRED Thresholds for Q2T1 and Q2T2 to 70 and 80, respectively
wrr-queue random-detect max-threshold 2 80 100
! Sets Max WRED Thresholds for Q2T1 and Q2T2 to 80 and 100, respectively

! This section maps the CoS values to the Queues/Thresholds
wrr-queue cos-map 1 1 1
! Maps CoS 1 (Scavenger/Bulk) to Q1 WRED Threshold 1
wrr-queue cos-map 1 2 0
! Maps CoS 0 (Best Effort) to Q1 WRED Threshold 2
wrr-queue cos-map 2 1 2 3
! Maps CoS 2 (Trans-Data & Mgmt) and CoS 3 (Call-Sig + Multimedia) to Q2T1
wrr-queue cos-map 2 2 6 7
! Maps CoS 6 (Routing) and CoS 7 (STP) to Q2 WRED Threshold 2
priority-queue cos-map 1 4 5
! Maps CoS 4 (TelePresence & Interactive-Video) and CoS 5 (VoIP) to the PQ
!

```

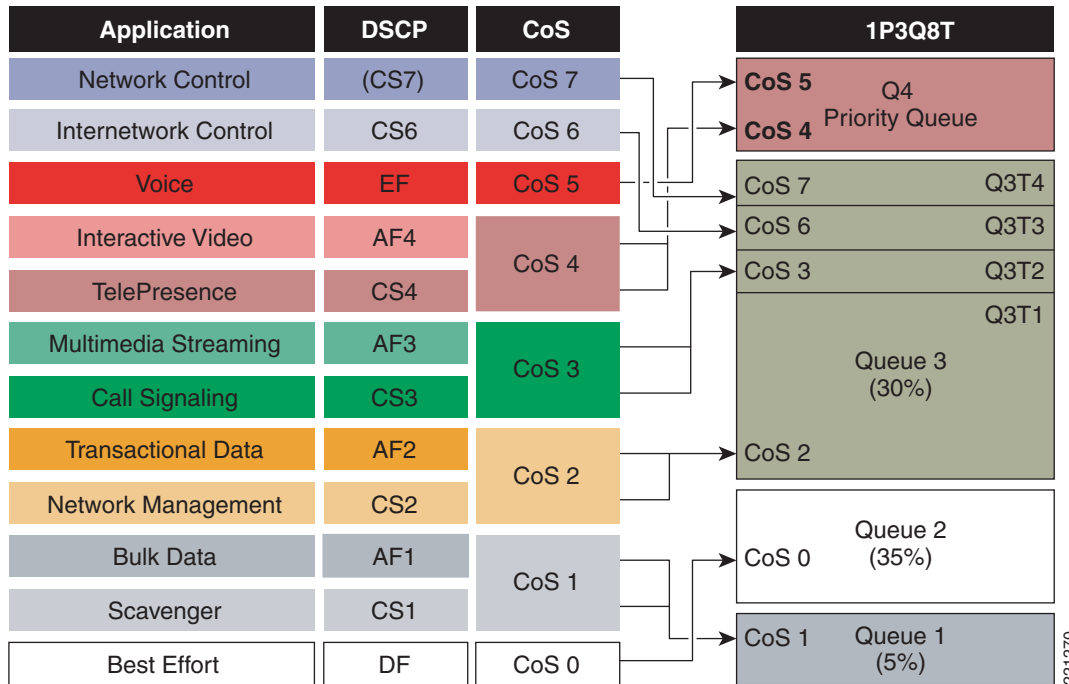
These configuration commands can be verified with the following command:

- **show queueing interface GigabitEthernet x/y**

Egress Queuing Design—1P3Q8T

As shown in [Figure 5-10](#), both the 6148A and 6748 linecards support a CoS-based egress queuing structure of 1P3Q8T, which uses WRED as a congestion avoidance mechanism. Under such a queuing structure, TelePresence traffic, along with VoIP, is recommended to be mapped to the strict-priority queue. Furthermore, because the egress queuing is CoS-based, Videoconferencing/Video Telephony (AF41) is also assigned to the strict-priority queue. CoS 1 (Scavenger/Bulk) can be constrained to a less than Best Effort queue: Q1. Q2 can then be dedicated for the default class (CoS 0). To minimize TCP global synchronization, WRED can be enabled on the non-realtime queues for congestion avoidance (technically, the congestion avoidance behavior is RED, as only one CoS weight is assigned to each queue). However, in Q3 the WRED thresholds can be set to give incremental preference to Network control traffic (CoS 7 and 6), followed by Call-Signaling traffic (CoS 3), and finally by Network Management traffic (CoS 2). The recommended egress 1P3Q8T queuing configuration for C6500 6148A and 6748 linecards is illustrated in [Figure 5-13](#).

Figure 5-13 Catalyst 6500 (1P3Q8T) Egress Queuing Recommendations for TelePresence Deployments



The configuration for the C6500 1P3Q8T egress queuing structure (for the 6548 linecards) illustrated in Figure 5-13 is as follows:

```
interface GigabitEthernet x/y

! This section sets the queue limits and bandwidth allocations
wrr-queue queue-limit 5 35 30
! Allocates 5% for Q1, 35% for Q2 and 30% for Q3
priority-queue queue-limit 30
! Allocates 30% for the Strict-Priority Queue (Q4)
wrr-queue bandwidth 5 35 30
! Sets the WRR weights for 5:35:30 (Q1:Q2:Q3) bandwidth servicing

! This section enables WRED on Q1, Q2 and Q3
wrr-queue random-detect 1
! Enables WRED on Q1
wrr-queue random-detect 2
! Enables WRED on Q2
wrr-queue random-detect 3
! Enables WRED on Q3

! This section sets Q1T1 WRED Thresholds to 80% (min) and 100% (max)
wrr-queue random-detect min-threshold 1 80 100 100 100 100 100 100 100
! Sets Min WRED Threshold for Q1T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 1 100 100 100 100 100 100 100 100
! Sets Max WRED Threshold for Q1T1 to 100% and all others to 100%

! This section sets Q2T1 WRED Thresholds to 80% (min) and 100% (max)
wrr-queue random-detect min-threshold 2 80 100 100 100 100 100 100 100
! Sets Min WRED Threshold for Q2T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 2 100 100 100 100 100 100 100 100
! Sets Max WRED Threshold for Q2T1 to 100% and all others to 100%
```

```

! This section sets Q3T1 to 60:70, Q3T2 to 70:80, Q3T3 to 80:90 and Q3T4 to 90:100
wrr-queue random-detect min-threshold 3 60 70 80 90 100 100 100 100
! Sets Min WRED Threshold for Q3T1 to 60%, Q3T2 to 70%, Q3T3 to 80%
! Q3T4 to 90%, and all others to 100%
wrr-queue random-detect max-threshold 3 70 80 90 100 100 100 100 100
! Sets Max WRED Threshold for Q3T1 to 70%, Q3T2 to 80%, Q3T3 to 90%
! and all others to 100%

! This section maps CoS values to egress Queues/Thresholds
wrr-queue cos-map 1 1 1
! Maps CoS 1 (Scavenger/Bulk) to Q1 WRED Threshold 1
wrr-queue cos-map 2 1 0
! Maps CoS 0 (Best Effort) to Q2 WRED Threshold 1
wrr-queue cos-map 3 1 2
! Maps CoS 2 (Net-Mgmt and Transactional Data) to Q3 WRED T1
wrr-queue cos-map 3 2 3
! Maps CoS 3 (Call-Signaling and Mission-Critical Data) to Q3 WRED T2
wrr-queue cos-map 3 3 6
! Maps CoS 6 (Routing) to Q3 WRED T3
wrr-queue cos-map 3 4 7
! Maps CoS 7 (Spanning Tree) to Q3 WRED T4
priority-queue cos-map 1 4 5
! Maps CoS 4 (TelePresence & Int-Video) and CoS 5 (VoIP) to the PQ
!

```

These configuration commands can be verified with the following command:

- **show queuing interface GigabitEthernet x/y**

Egress Queuing Design—1P7Q8T

As shown in [Figure 5-10](#), the 6704-10GE linecards support a CoS-based egress queuing structure of 1P7Q8T, which uses WRED as a congestion avoidance mechanism. Such a queuing structure allows for a dedicated queue for each CoS value.

Normally, TelePresence traffic (represented by CoS 4, along with Videoconferencing) would be mapped to the strict priority queue in campus queuing models; however, as a dedicated hardware queue exists on this platform, an exception could be made to this rule. As long as both TelePresence and Videoconferencing traffic are bounded by Call-Admission Control (CAC) mechanisms, these can adequately be provisioned in a dedicated, non-priority queue and still meet the tight service levels required by TelePresence. Remember, in the campus, latency and jitter are so minimal that they are considered negligible (typically in the microseconds on a per-switch basis). However loss is the most important service level consideration within the campus, since packet loss can easily occur within a campus due to instantaneous buffer congestion. The key is to ensure that switching buffers are never overrun during such instantaneous bursts. This can be achieved either by assigning TelePresence traffic to a (sufficiently deep) priority hardware queue or by assigning TelePresence to a non-priority queue, which has a limit to how much traffic can be admitted to the queue (via a CAC mechanism).

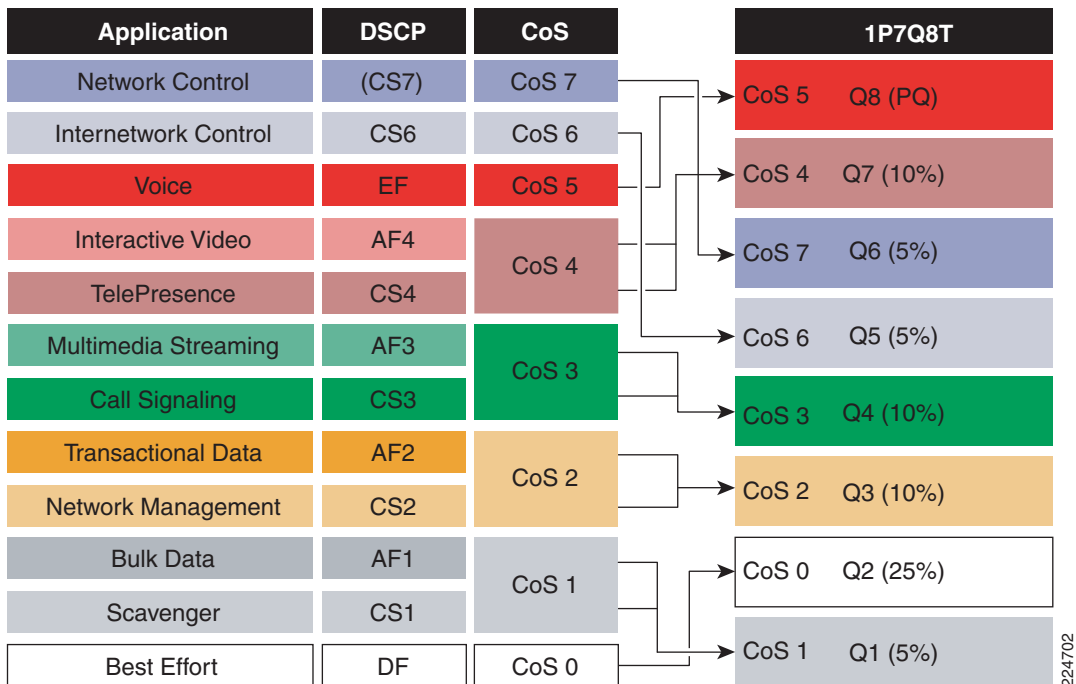
In this example, we present the option of assigning TelePresence (along with Videoconferencing) traffic to a dedicated non-priority queue (Q7), whereas VoIP continues to be mapped to the strict-priority queue (Q8). All other CoS-based traffic classes can be mapped to their dedicated queues.

To minimize TCP global synchronization, WRED can be enabled on the non-realtime queues for congestion avoidance (technically, the congestion avoidance behavior is RED, as only one CoS weight is assigned to each queue). This can be done by setting the minimum threshold for WRED to 80% and the maximum threshold to 100% (the tail of the queue) on a per-queue basis. However, there is no value—and only harm—from enabling WRED on the dedicated, non-priority queue for

TelePresence/Videoconferencing (Q7), as you never want to aggressively drop TelePresence traffic for any reason. As such, WRED can be disabled by setting the minimum and maximum WRED threshold to 100% (the tail of the queue).

The recommended egress 1P7Q8T queuing configuration for C6500 6704-10GE linecards is illustrated in Figure 5-14.

Figure 5-14 Catalyst 6500 (1P7Q8T) Egress Queuing Recommendations for TelePresence Deployments



The configuration for the C6500 1P7Q8T egress queuing structure (for the 6704-10GE linecards) illustrated in Figure 5-14 is as follows:

```

!
interface range TenGigabitEthernet4/1 - 4
  wrr-queue queue-limit 5 25 10 10 5 5 10
    ! Allocates 5% to Q1, 25% to Q2, 10% to Q3, 10% to Q4,
    ! Allocates 5% to Q5, 5% to Q6 and 10% to Q7
  wrr-queue bandwidth 5 25 10 10 5 5 10
    ! Sets the WRR weights for 5:25:10:10:5:5:10 (Q1 through Q7)
  priority-queue queue-limit 30
    ! Limits PQ to 30%

  wrr-queue random-detect 1
    ! Enables WRED on Q1
  wrr-queue random-detect 2
    ! Enables WRED on Q2
  wrr-queue random-detect 3
    ! Enables WRED on Q3
  wrr-queue random-detect 4
    ! Enables WRED on Q4
  wrr-queue random-detect 5
    ! Enables WRED on Q5
  wrr-queue random-detect 6
    ! Enables WRED on Q6
  wrr-queue random-detect 7

```

```

! Enables WRED on Q7

wrr-queue random-detect min-threshold 1 80 100 100 100 100 100 100
! Sets Min WRED Threshold for Q1T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 1 100 100 100 100 100 100 100
! Sets Max WRED Threshold for Q1T1 to 100% and all others to 100%

wrr-queue random-detect min-threshold 2 80 100 100 100 100 100 100
! Sets Min WRED Threshold for Q2T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 2 100 100 100 100 100 100 100
! Sets Max WRED Threshold for Q2T1 to 100% and all others to 100%

wrr-queue random-detect min-threshold 3 80 100 100 100 100 100 100
! Sets Min WRED Threshold for Q3T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 3 100 100 100 100 100 100 100
! Sets Max WRED Threshold for Q3T1 to 100% and all others to 100%

wrr-queue random-detect min-threshold 4 80 100 100 100 100 100 100
! Sets Min WRED Threshold for Q4T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 4 100 100 100 100 100 100 100
! Sets Max WRED Threshold for Q4T1 to 100% and all others to 100%

wrr-queue random-detect min-threshold 5 80 100 100 100 100 100 100
! Sets Min WRED Threshold for Q5T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 5 100 100 100 100 100 100 100
! Sets Max WRED Threshold for Q5T1 to 100% and all others to 100%

wrr-queue random-detect min-threshold 6 80 100 100 100 100 100 100
! Sets Min WRED Threshold for Q6T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 6 100 100 100 100 100 100 100
! Sets Max WRED Threshold for Q6T1 to 100% and all others to 100%

wrr-queue random-detect min-threshold 7 100 100 100 100 100 100 100
! Sets all Min WRED Thresholds for Q7 to 100 (Disables WRED on Q7)
wrr-queue random-detect max-threshold 7 100 100 100 100 100 100 100
! Sets all Max WRED Thresholds for Q7 to 100 (Disables WRED)
! WRED is disabled on TelePresence Queue (Q7); Tail-drop only

wrr-queue cos-map 1 1 1
! Maps Scavenger/Bulk to Q1T1
wrr-queue cos-map 2 1 0
! Maps Best Effort to Q2T1
wrr-queue cos-map 3 1 2
! Maps Trans-Data & Net-Mgmt to Q3T1
wrr-queue cos-map 4 1 3
! Maps MM-Streaming & Call-Signaling to Q4T1
wrr-queue cos-map 5 1 6
! Maps IP Routing to Q5T1
wrr-queue cos-map 6 1 7
! Maps Spanning Tree to Q6T1
wrr-queue cos-map 7 1 4
! Maps TelePresence to Q7T1 (WRED is Disabled)
priority-queue cos-map 1 5
! Maps Voice to the PQ (Q8)
mls qos trust dscp
! Sets interface to trust DSCP
!

```

**Note**

Due to WRED defaults and/or settings, the order of entering WRED min/max threshold commands may need to be reversed for some queues. The commands have been shown in the order presented to emphasize and simplify the logic of the configuration, rather than to appease the IOS parser.

These configuration commands can be verified with the following command:

- **show queueing interface GigabitEthernet x/y**

Egress Queuing Design—1P7Q4T (DSCP-to-Queue)

As shown in [Figure 5-10](#), the 6708-10GE linecards support a CoS- or DSCP-based egress queuing structure of 1P7Q4T, which uses WRED as a congestion avoidance mechanism.

The increased granularity in queue-mapping that is presented by the support of DSCP-to-Queue mapping allows—for the first time—TelePresence traffic to be separated from Videoconferencing traffic. As such, TelePresence can be assigned to its own dedicated, non-priority queue (Q7), separate from the dedicated, non-priority queue for Videoconferencing (Q5). Again, if a non-priority queue is to be used to service TelePresence traffic, then it is essential that a Call Admission Control mechanism be used to limit TelePresence traffic, inline with the bandwidth assigned to this queue.

VoIP (DSCP EF) will continue to be mapped to the strict-priority queue (Q8).

Queue 6 can be used to service control traffic, including (in respective order):

- Network Control traffic (DSCP CS7)
- Internetwork Control traffic (DSCP 6)
- Call Signaling traffic (DSCP CS3)
- Network Management traffic (CS2)

This respective order can be enforced by leveraging the 4 drop thresholds that can be assigned per non-priority queue. Specifically, CS2 can be mapped to Q6T1, CS3 can be mapped to Q6T2, CS6 can be mapped to Q6T3, and CS7 can be mapped to Q6T4 (the tail of the queue).

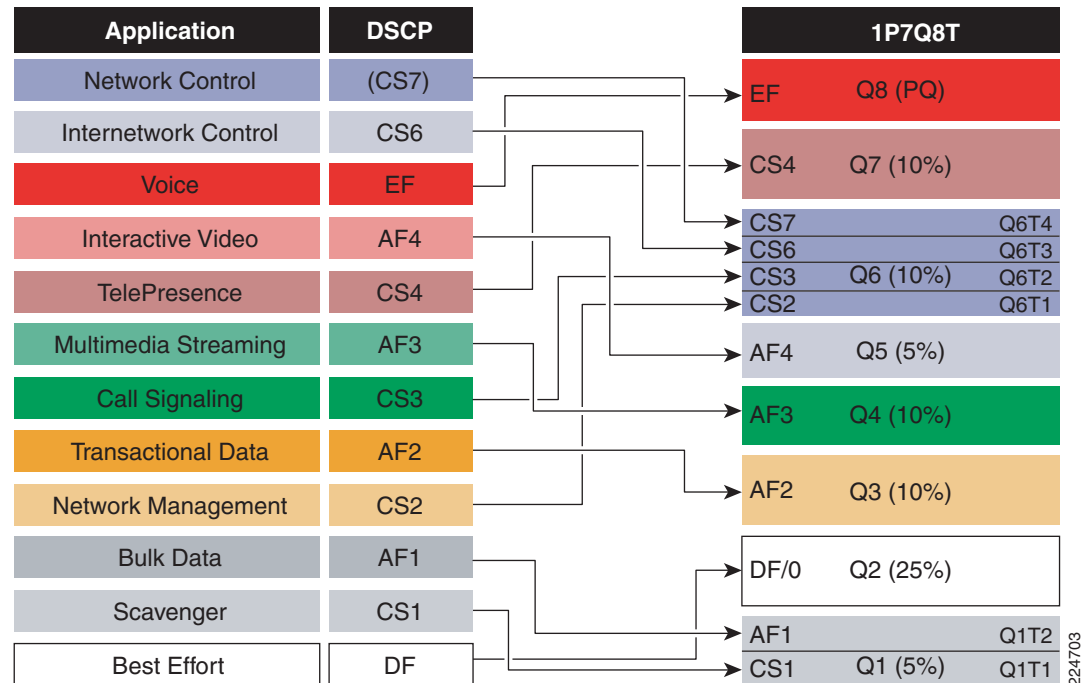
Multimedia Streaming (AF3) traffic can be assigned to a dedicated, non-priority queue (Q4), as can be Transactional Data (AF2) traffic (to Q3). This leaves Q2 to serve as a default/Best Effort queue and Q1 to serve as a “less than Best Effort” queue to service Bulk Data (AF1) and Scavenger (CS1) traffic, respectively. Again, the respective priority of Bulk Data over Scavenger traffic in Q1 can be enforced by leveraging the configurable WRED drop thresholds for the queue, such that Scavenger is mapped to Q1T1 and Bulk Data to Q1T2.

To minimize TCP global synchronization, WRED can be enabled on the non-realtime queues for congestion avoidance. Additionally, as this platform supports DSCP-to-Queue/Threshold mappings, the Assured Forwarding Per-Hop Behavior (AF PHB) can be fully-implemented on all AF traffic classes by enabling DSCP-based WRED (such that AFx3 drops before AFx2, which-respectively-is dropped before AFx1).

However, as noted before, there is no value—and only harm—from enabling WRED on the dedicated, non-priority queue for TelePresence (Q7), as you never want to aggressively drop TelePresence traffic for any reason. As such, WRED can be disabled by setting the minimum and maximum WRED threshold to 100% (the tail of the queue).

The recommended egress 1P7Q4T (DSCP-to-Queue) queuing configuration for C6500 6708-10GE linecards is illustrated in [Figure 5-15](#).

Figure 5-15 Catalyst 6500 (1P7Q4T—DSCP-to-Queue) Egress Queuing Recommendations for TelePresence Deployments



224703

The configuration for the C6500 1P7Q4T (DSCP-to-Queue) egress queuing structure (for the 6708-10GE linecards) illustrated in Figure 5-15 is as follows:

```

!
interface range TenGigabitEthernet4/1 - 8
  wrr-queue queue-limit 5 25 10 10 10 10 10
    ! Allocates 5% to Q1, 25% to Q2, 10% to Q3, 10% to Q4,
    ! Allocates 5% to Q5, 5% to Q6 and 10% to Q7
  wrr-queue bandwidth 5 25 10 10 10 10 10
    ! Sets the WRR weights for 5:25:10:10:10:10:10 (Q1 through Q7)
  priority-queue queue-limit 20
    ! Limits PQ to 20%

  wrr-queue random-detect 1
    ! Enables WRED on Q1
  wrr-queue random-detect 2
    ! Enables WRED on Q2
  wrr-queue random-detect 3
    ! Enables WRED on Q3
  wrr-queue random-detect 4
    ! Enables WRED on Q4
  wrr-queue random-detect 5
    ! Enables WRED on Q5
  wrr-queue random-detect 6
    ! Enables WRED on Q6
  wrr-queue random-detect 7
    ! Enables WRED on Q7

  wrr-queue random-detect min-threshold 1 60 70 80 90 100 100 100
    ! Sets Min WRED Thresholds for Q1T1 to 60%, Q1T2 to 70%, Q1T3 to 80%, Q1T4 to 90%
  wrr-queue random-detect max-threshold 1 70 80 90 100 100 100 100
    ! Sets Max WRED Thresholds for Q1T1 to 70%, Q1T2 to 80%, Q1T3 to 90%, Q1T4 to 100%

  wrr-queue random-detect min-threshold 2 80 100 100 100 100 100 100

```

```

! Sets Min WRED Threshold for Q2T1 to 80% and all others to 100%
wrr-queue random-detect max-threshold 2 100 100 100 100 100 100 100
! Sets all Max WRED Thresholds for Q2 to 100%

wrr-queue random-detect min-threshold 3 70 80 90 100 100 100 100
! Sets Min WRED Thresholds for Q3T1 to 70%, Q3T2 to 80%, Q3T3 to 90%, Q3T4 to 100%
wrr-queue random-detect max-threshold 3 80 90 100 100 100 100 100
! Sets Max WRED Thresholds for Q3T1 to 80%, Q3T2 to 90%, Q3T3 & Q3T4 to 100%

wrr-queue random-detect min-threshold 4 70 80 90 100 100 100 100
! Sets Min WRED Thresholds for Q4T1 to 70%, Q4T2 to 80%, Q4T3 to 90%, Q4T4 to 100%
wrr-queue random-detect max-threshold 4 80 90 100 100 100 100 100
! Sets Max WRED Thresholds for Q4T1 to 80%, Q4T2 to 90%, Q4T3 & Q4T4 to 100%

wrr-queue random-detect min-threshold 5 70 80 90 100 100 100 100
! Sets Min WRED Thresholds for Q5T1 to 70%, Q5T2 to 80%, Q5T3 to 90%, Q5T4 to 100%
wrr-queue random-detect max-threshold 5 80 90 100 100 100 100 100
! Sets Max WRED Thresholds for Q5T1 to 80%, Q5T2 to 90%, Q5T3 & Q5T4 to 100%

wrr-queue random-detect min-threshold 6 60 70 80 90 100 100 100
! Sets Min WRED Thresholds for Q6T1 to 60%, Q6T2 to 70%, Q6T3 to 80%, Q6T4 to 90%
wrr-queue random-detect max-threshold 6 70 80 90 100 100 100 100
! Sets Max WRED Thresholds for Q6T1 to 70%, Q6T2 to 80%, Q6T3 to 90%, Q6T4 to 100%

wrr-queue random-detect min-threshold 7 100 100 100 100 100 100 100
! Sets all Min WRED Thresholds for Q7 to 100 (Disables WRED on Q7)
wrr-queue random-detect max-threshold 7 100 100 100 100 100 100 100
! Sets all Max WRED Thresholds for Q7 to 100 (Disables WRED)
! WRED is disabled on TelePresence Queue (Q7); Tail-drop only

mls qos trust dscp
! Sets interface to trust DSCP
mls qos queue-mode mode-dscp
! Enables DSCP-to-Queue mapping mode

wrr-queue dscp-map 1 1 8
! Maps Scavenger (CS1) to Q1T1
wrr-queue dscp-map 1 2 14
! Maps Bulk Data (AF13) to Q1T2 - AF13 PHB
wrr-queue dscp-map 1 3 12
! Maps Bulk Data (AF12) to Q1T3 - AF12 PHB
wrr-queue dscp-map 1 4 10
! Maps Bulk Data (AF11) to Q1T4 - AF11 PHB

wrr-queue dscp-map 2 1 0
! Maps Best Effort to Q2T1

wrr-queue dscp-map 3 1 22
! Maps Transactional Data (AF23) to Q3T1 - AF23 PHB
wrr-queue dscp-map 3 2 20
! Maps Transactional Data (AF22) to Q3T2 - AF22 PHB
wrr-queue dscp-map 3 3 18
! Maps Transactional Data (AF21) to Q3T3 - AF21 PHB

wrr-queue dscp-map 4 1 30
! Maps Multimedia Streaming (AF33) to Q4T1 - AF33 PHB
wrr-queue dscp-map 4 2 28
! Maps Multimedia Streaming (AF32) to Q4T2 - AF32 PHB
wrr-queue dscp-map 4 3 26
! Maps Multimedia Streaming (AF31) to Q4T3 - AF31 PHB

wrr-queue dscp-map 5 1 38
! Maps Multimedia Conferencing (AF43) to Q5T1 - AF43 PHB

```

```
wrr-queue dscp-map 5 2 36
! Maps Multimedia Conferencing (AF42) to Q5T2 - AF42 PHB
wrr-queue dscp-map 5 3 34
! Maps Multimedia Conferencing (AF41) to Q5T3 - AF41 PHB

wrr-queue dscp-map 6 1 16
! Maps Net-Mgmt (CS2) to Q6T1
wrr-queue dscp-map 6 2 24
! Maps Call-Signaling (CS3) to Q6T2
wrr-queue dscp-map 6 3 48
! Maps IP Routing (CS6) to Q6T3
wrr-queue dscp-map 6 4 56
! Maps Spanning-Tree (CS7) to Q6T4

wrr-queue dscp-map 7 1 32
! Maps TelePresence (CS4) to Q7T1

priority-queue dscp-map 1 46
! Maps Voice (EF) to the PQ (Q8)
!
```

**Note**

Due to WRED defaults and/or settings, the order of entering WRED min/max threshold commands may need to be reversed for some queues. The commands have been shown in the order presented to emphasize and simplify the logic of the configuration, rather than to appease the IOS parser.

These configuration commands can be verified with the following command:

- **show queueing interface GigabitEthernet x/y**



CHAPTER 6

Branch QoS Design for TelePresence

TelePresence Branch QoS Design Overview

The primary business advantages of TelePresence systems include:

- Reduced travel time and expense
- Improved collaboration and productivity
- Improved quality of work/life (due to reduced travel)
- The green advantage of a reduced carbon footprint

However, these business advantages are not fully realized if TelePresence systems are connected solely via an Intra-Campus Deployment Model (as illustrated in [Figure 3-1](#)); rather, gaining these advantages requires TelePresence systems to be deployed over wide area networks, whether these are private WANs or Virtual Private Networks.

WANs or VPNs may be used to interconnect large campuses to each other or may be used to connect one or more large campuses with smaller branch offices (as illustrated in [Figure 3-3](#)). To simplify these permutations, we refer to all TelePresence connections over a wide area as Branch Places-in-the-Network (PINs).

Branch PINs serve as boundary points between local area and wide area networks and, as such, these are often the most bottlenecked PINs and therefore have the most critical QoS requirements within the network infrastructure. To help select the best policies to be used at these critical PINs, it is beneficial to review some important considerations, which we discuss next.

LLQ versus CBWFQ Considerations

Probably the most controversial decision relating to TelePresence deployments is whether to provision TelePresence traffic over the WAN/VPN in a strict-priority Low-Latency Queue (LLQ) or in a dedicated bandwidth-guaranteed Class-Based Weighted-Fair Queue (CBWFQ).

In campus networks, placing TelePresence in the strict-priority hardware queues yielded superior results during testing, especially in terms of protection against packet loss during momentary periods of congestion, which occur regularly in campus networks even under normal operating conditions. Additionally, placing TelePresence traffic in these strict-priority queues does not involve any incremental or ongoing monetary expense (beyond initial configuration), as this potential for strict-priority servicing already exists within the campus network infrastructure and the exercise simply becomes a matter of re-configuring existing queuing structures to enable strict-priority queuing for TelePresence.

Therefore, the decision to service TelePresence with strict-priority queues within the campus is relatively straightforward.

However, the corresponding decision becomes more complicated over the WAN/VPN due to three main considerations:

- The cost of subscribing to realtime SP services
- The “33% LLQ Rule”
- The potential effect of TelePresence on VoIP

Let us look at each of these considerations in turn. The first and foremost consideration is the ongoing cost of subscribing to realtime services from a service provider. Service providers generally charge enterprise customers premium rates for the amount of traffic they want serviced within a realtime class. At times, these additional premiums may make it cost prohibitive to provision TelePresence traffic within a realtime SP class. At the very least, such expensive premiums could diminish the overall business cost savings that TelePresence can provide an enterprise (versus employee travel expenses).

The second consideration is the potential impact of the “33% LLQ Rule” (referenced in [Chapter 4, “Quality of Service Design for TelePresence”](#) in the section [Queuing TelePresence](#)). At times, administrators cannot provision adequate amounts of bandwidth for TelePresence and remain within this conservative design recommendation. This is generally the case when dealing with (45 Mbps) T3/DS3 links. According to the “33% LLQ Rule,” no more than 15 Mbps of traffic of such a link should be assigned for strict-priority servicing. However, if a network administrator already has VoIP provisioned (quite properly) in an LLQ on such a link, and is looking to also provision TelePresence with strict-priority servicing, then they have a decision to make. For example, if they wish to deploy a CTS-3000 at 1080p-Best (requiring 15 Mbps just for TelePresence), then they either have to upgrade the link’s bandwidth capacity (which is often cost-prohibitive, as generally the next tier of bandwidth is OC3) or they violate this design rule to accommodate all of their realtime traffic.

At this point, it bears repeating that the “33% LLQ Rule” is a conservative design recommendation, with the intent of reducing the variance in application response times of non-realtime applications during periods that the realtime classes are being utilized at maximum capacity. This is an exceptionally relevant concern when dealing with a high-bandwidth realtime application such as TelePresence.

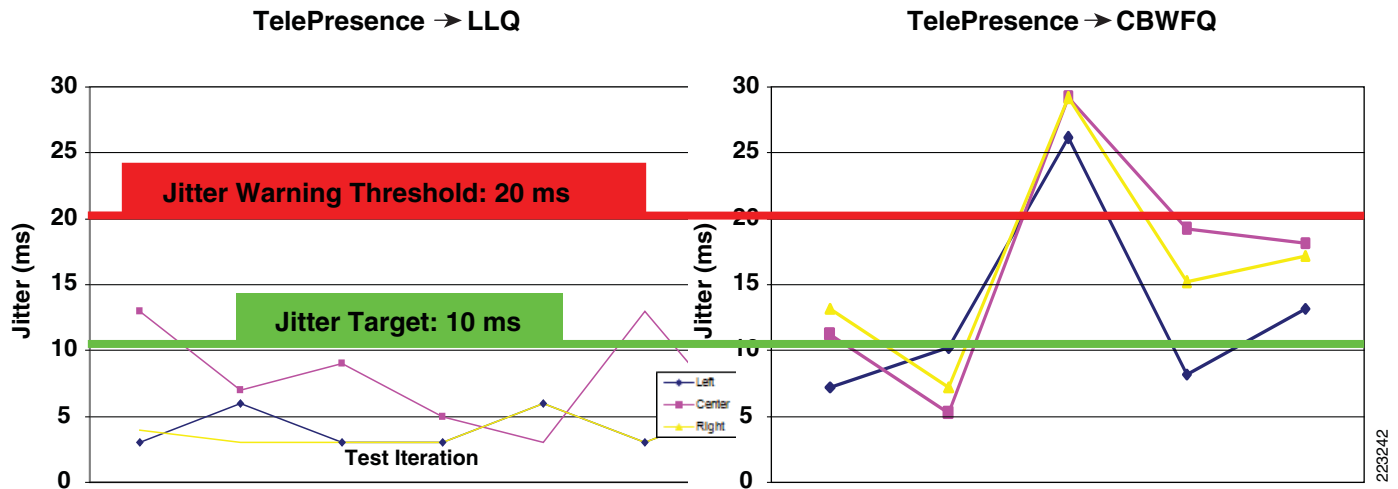
For example, let us reconsider a (45 Mbps) T3/DS3 link configured to support two separate CTS-3000 calls, both configured to transmit at full 1080p-Best resolution. Each such call requires 15 Mbps of realtime traffic. Prior to TelePresence calls being placed, non-realtime applications would have access to 100% of the bandwidth (to simplify the example, we are assuming there are no other realtime applications, such as VoIP, on this link). However, once these TelePresence calls are established, realtime TelePresence calls would suddenly dominate more than 66% of the link and all non-realtime applications would just as suddenly be contending for less than 33% of the link. TCP windowing for many of these non-realtime applications would begin slow-starting, resulting in many data applications hanging, timing out, or becoming stuck in a non-responsive state. Such network behavior, changing from one minute to the next, generally translates into users calling the IT help desk complaining about the network (which happens to be functioning properly, albeit in a poorly-configured manner).

That being said, it bears repeating that the “33% LLQ Rule” rule is not to be viewed as a mandate, but is simply a best practice design recommendation. There may be cases where specific business objectives cannot be met while holding to this recommendation. In such cases, enterprises must provision according to their detailed requirements and constraints. However, it is important to recognize the tradeoffs involved with over-provisioning realtime traffic classes in conjunction with the negative performance impact this has on non-realtime-application response times.

Quite naturally then, to make such provisioning decisions, a network administrator might wonder about the tradeoffs involved in TelePresence application performance when TelePresence is placed in a LLQ versus a CBWFQ. In such a comparison, the most sensitive service level attribute is jitter, as both

policies can be configured to completely prevent packet loss. As one might suspect, provisioning TelePresence in a LLQ results in lower peak-to-peak jitter values, as compared to provisioning TelePresence in a CBWFQ, which is shown in Figure 6-1.

Figure 6-1 Jitter Comparisons Between LLQ and CBWFQ WAN Edge Queuing Policies



Cisco Enterprise System Engineering testing showed that a 12-class RFC 4594-based QoS policy on a fully-congested T3 link—with TelePresence being serviced in a LLQ—yielded between 3-13 ms of peak-to-peak jitter to TelePresence; whereas an identical test—but with TelePresence being serviced in a CBWFQ—yielded between 5-29 ms of peak-to-peak jitter to TelePresence. As detailed in [Chapter 4, “Quality of Service Design for TelePresence,”](#) TelePresence has a one-way peak-to-peak jitter target of 10 ms and a warning threshold of 20 ms of peak-to-peak jitter, which if exceeded over an extended period can generate warning messages on the screen. During some of the CBWFQ tests, this warning message was observed on the screen, indicating that network congestion was affecting the TelePresence call quality.



Note

These tests were performed using TelePresence codec software version 1.1.0 (256D). Newer versions of CTS software have superior traffic smoothing capabilities as well as deeper de-jitter buffering, both of which amount to less overall sensitivity of TelePresence to jitter. Therefore the advantage of LLQ over CBWFQ queuing policies is less with newer versions of CTS software.

Therefore, while a moderate performance advantage to TelePresence can be observed when it is provisioned in a LLQ versus a CBWFQ, the advantage is not so great as to preclude recommending provisioning TelePresence in a CBWFQ when it is not viable to be provisioned with a LLQ. In other words, from a purely technical standpoint, **the best performance levels for TelePresence can be achieved when it is provisioned in a LLQ.** However, when other factors (such as additional ongoing costs or over-provisioning constraints for realtime bandwidth, etc.) need to be taken into account and render provisioning TelePresence in an LLQ unviable, then **the next best levels of service can be achieved by provisioning TelePresence in a dedicated CBWFQ.**

The third main consideration of whether to use LLQ or CBWFQ is the potential effect of TelePresence traffic on VoIP traffic if both are to be serviced in a strict-priority queue. To better understand how these realtime applications can be provisioned with strict-priority servicing and protected from interfering with each other, we must take a closer look at Cisco’s IOS LLQ/CBWFQ mechanisms. To do so, let us consider a simple LLQ/CBWFQ policy.

Example 6-1 Simple LLQ/CBWFQ Policy

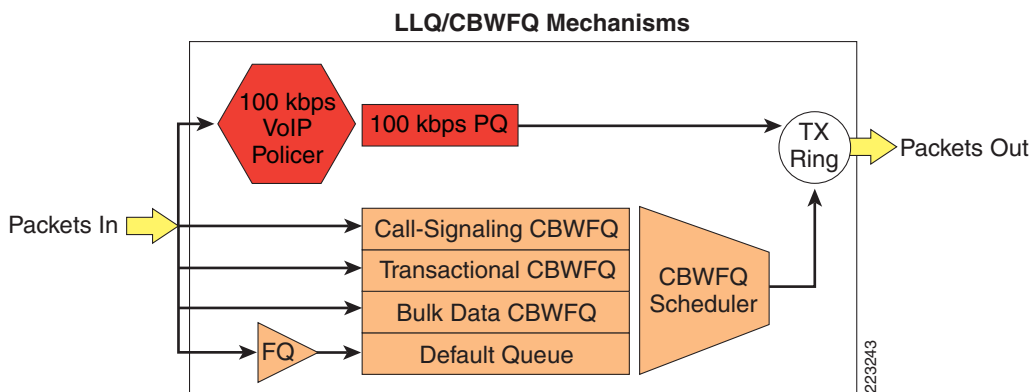
```

policy-map WAN-EDGE
  class VOIP
    priority 100
  class CALL-SIGNALING
    bandwidth percent 5
  class TRANSACTIONAL
    bandwidth percent 20
  class BULK
    bandwidth percent 10
  class class-default
    fair-queue

```

The underlying mechanisms for this LLQ/CBWFQ policy are graphically represented in [Figure 6-2](#).

Figure 6-2 Cisco IOS LLQ/CBWFQ Mechanisms—Part 1

**Note**

For the sake of simplicity, some Layer 2 subsystems (including Link Fragmentation and Interleaving) have been omitted from [Figure 6-2](#), as these mechanisms simply are not relevant at the link speeds required by TelePresence.

In [Figure 6-2](#), we see a router interface that has been configured with a 5-class LLQ/CBWFQ policy, with VoIP assigned to a 100 kbps LLQ and additional three explicit CBWFQs defined for Call-Signaling, Transactional Data, and Bulk Data respectively, as well as a default queue that has a Fair-Queuing pre-sorter assigned to it. There are two additional underlying mechanisms that may not be obvious from the configuration, but are shown in [Figure 6-2](#):

- An implicit policer attached to the LLQ
- A final output buffer called the Tx-Ring

Let us first take a look at the implicit policer attached to the LLQ. The threat posed by any strict priority-scheduling algorithm is that it could completely starve lower priority traffic. To prevent this, the LLQ mechanism has a built-in policer. This policer (like the queuing algorithm itself) engages only when the interface is experiencing congestion. Therefore, it is important to provision the priority classes properly. In this example, if more than 100 kbps of VoIP traffic was offered to the interface and the interface was congested, the excess VoIP traffic would be discarded by the implicit policer. However, traffic admitted by the policer gains access to the strict priority queue and is handed off to the Tx-Ring ahead of all other CBWFQ traffic.

The Tx-Ring is a final output buffer that serves the purpose of always having packets ready to be placed onto the wire so that link utilization can be driven to 100%. It is actually a full Tx-Ring that signals the Cisco IOS software to indicate that an interface is experiencing congestion and as such the LLQ/CBWFQ algorithms need to be engaged.

Now, let us consider the case of servicing not just VoIP with strict-priority queuing, but also TelePresence.

**Note**

For the sake of example and illustration simplicity, let us assume TelePresence only requires 400 kbps of traffic for these next two examples only.

Two options exist to the network administrator. The first is to admit both VoIP and TelePresence to the same LLQ. Thus our example policy becomes:

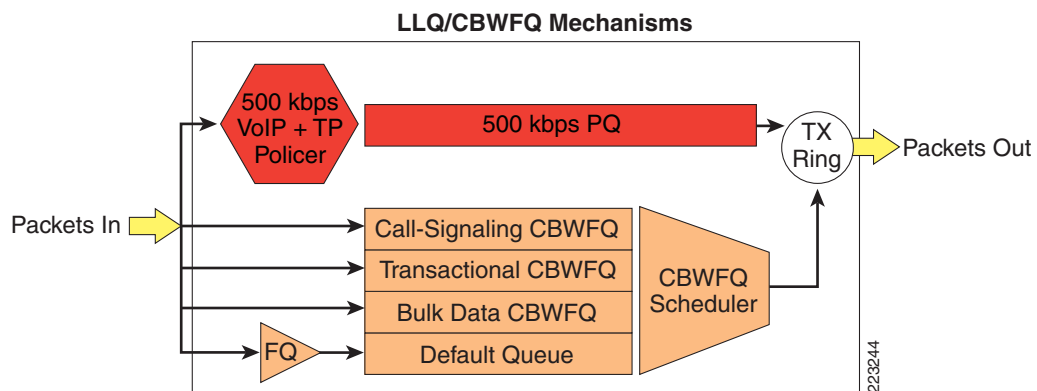
Example 6-2 VoIP and TelePresence in a Single LLQ Policy

```
class-map match-any REALTIME
  match dscp ef          ! Matches VoIP
  match dscp cs4        ! Matches TelePresence
  ...

policy-map WAN-EDGE
  class REALTIME
    priority 500         ! 100 kbps for VoIP + 400 kbps for TelePresence
  class CALL-SIGNALING
    bandwidth percent 5
  class TRANSACTIONAL
    bandwidth percent 20
  class BULK
    bandwidth percent 10
  class class-default
    fair-queue
```

The corresponding IOS mechanisms for [Example 6-2](#) are illustrated in [Figure 6-3](#).

Figure 6-3 Cisco IOS LLQ/CBWFQ Mechanisms—Part 2



In [Figure 6-3](#), we can see that not only has the LLQ been expanded in size (to 500 kbps), but also the implicit policer (for the combined VoIP and TelePresence class) has been increased to 500 kbps. Such a policy continues to protect VoIP from data as well as TelePresence from data. However, this policy does

potentially allow TelePresence to interfere with VoIP. This is because traffic offered to the LLQ class is serviced on a first-come, first-serve basis. Therefore, should TelePresence traffic suddenly burst, then it is possible—even likely—that VoIP traffic would be dropped.

At this point, we can realize another benefit of the implicit policer for the LLQ: not only does this mechanism protect non-realtime queues from bandwidth-starvation, but also it allows for Time-Division Multiplexing (TDM) of the LLQ. TDM of the LLQ allows for the configuration and servicing of “multiple” LLQs, while abstracting the fact that there is only a single LLQ “under-the-hood,” so to speak. Pertinent to our example, by configuring two LLQs, not only are VoIP and TelePresence protected from data applications, but VoIP and TelePresence are also protected from interfering with each other.

Let us take a look at our final policy example to cover this point. In [Example 6-3](#), a dual-LLQ design is used, one each for VoIP and TelePresence.

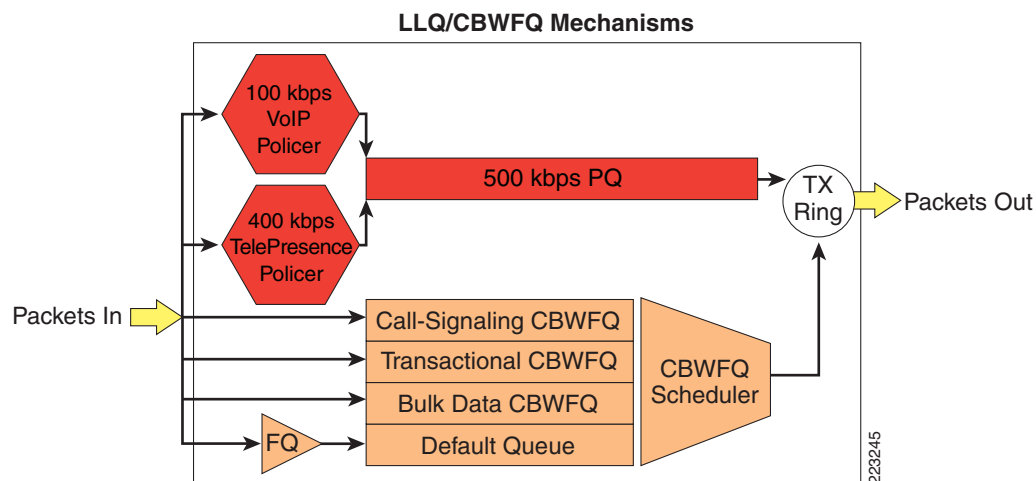
Example 6-3 VoIP and TelePresence in a Dual-LLQ Policy

```
class-map match-all VOIP
  match dscp ef                ! Matches VoIP
class-map match-all TELEPRESENCE
  match dscp cs4              ! Matches TelePresence
...

policy-map WAN-EDGE
  class VOIP
    priority 100              ! 100 kbps LLQ for VoIP
  class TELEPRESENCE
    priority 400              ! 400 kbps LLQ for TelePresence
  class CALL-SIGNALING
    bandwidth percent 5
  class TRANSACTIONAL
    bandwidth percent 20
  class BULK
    bandwidth percent 10
  class class-default
    fair-queue
```

The corresponding IOS mechanisms for [Example 6-3](#) are illustrated in [Figure 6-4](#).

Figure 6-4 Cisco IOS LLQ/CBWFQ Mechanisms—Part 3



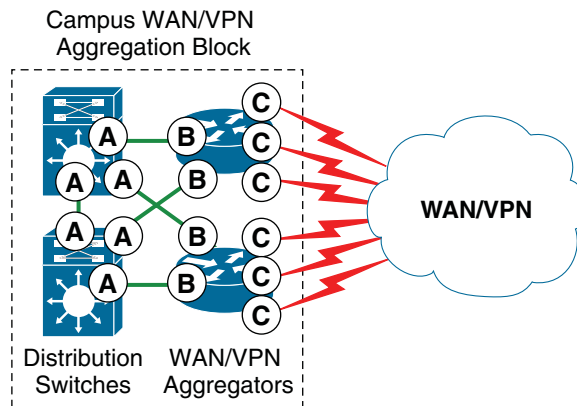
In Figure 6-4, we see that two separate implicit policers have been provisioned, one each for the VoIP class (to 100 kbps) and another for the TelePresence class (to 400 kbps), yet there remains only a single strict-priority queue, which is provisioned to the sum of all LLQ classes, in this case to 500 kbps (100 kbps + 400 kbps). Traffic offered to either LLQ class is serviced on a first-come, first-serve basis until the implicit policer for each specific class has been invoked. For example, if TelePresence attempts to burst beyond a 400 kbps rate (remember, this rate has been reduced in order to simplify this example, both textually and also graphically), then it is dropped.

Therefore, to sum up this final consideration regarding whether or not to use LLQ for TelePresence: **if strict priority servicing for TelePresence is desired and viable, and if another realtime class (such as VoIP) has already been configured with a LLQ, then a dual-LLQ design would be recommended** in order to protect VoIP and TelePresence from interfering with each other.

Campus WAN/VPN Block Considerations

Typically the first step in connecting a branch to a campus is to build out a WAN/VPN aggregation block at the main campus site. An example enterprise campus WAN/VPN aggregation block is illustrated in Figure 6-5.

Figure 6-5 Enterprise Campus WAN/VPN Aggregation Block QoS Design Recommendations for TelePresence



- A Interswitch Link Policies:**
Trust DSCP
+ Queuing (CoS 4 and 5 → PQ)
+ Queuing (CoS 3 → Non-PQ)
- B Router LAN Edge Policies:**
Trust DSCP (default)
+ LLQ for VoIP (EF)
+ LLQ for TelePresence (CS4)
+ CBWFQ for Call-Signaling (CS3)
- C WAN/VPN Edge QoS Policies:**
Trust for VoIP (EF)
+ LLQ or CBWFQ for TelePresence (CS4)
+ CBWFQ for Call-Signaling (CS3)

223246

Interswitch links are shown in [Figure 6-5](#) as points labeled A. While technically-speaking some of these links are interconnecting switches to routers, however their role and configuration are the same as the interswitch links described and defined in [Chapter 5, “Campus QoS Design for TelePresence.”](#) To be completely technically accurate, we could refer to these links as LAN-to-LAN non-edge links, but this term becomes a bit wordy and unwieldy, and as such we continue to use the simpler term interswitch links, but with a broadened meaning to include these switch-to-router links as well.

**Note**

As noted above, the term used here as interswitch links in this context refers to LAN-to-LAN non-edge links, not (necessarily) trunked links encapsulated with Cisco InterSwitch Link (ISL) trunking protocol.

As previously detailed (in [Chapter 5, “Campus QoS Design for TelePresence”](#)), all interswitch links should be configured to trust DSCP and perform hardware queuing, such that CoS 4 (TelePresence) and CoS 5 (VoIP) are assigned to the strict priority hardware queue and CoS 3 (Call-Signaling) is assigned to a non-priority queue within the platform/linecard’s 1PxDyT queuing structure.

Next, it would be recommended to enable LLQ/CBWFQ (or hardware queuing policies, if supported) on the router’s LAN edges, labeled as points B in [Figure 6-5](#). This is recommended when the levels of WAN/VPN aggregation may make it theoretically possible to oversubscribe these WAN-to-LAN links. For example, if the WAN/VPN aggregation router was homing seven individual OC3 circuits (totaling 7 * 155 Mbps or 1.085 Gbps), but connecting to the distribution switches via GigabitEthernet links, then the potential for oversubscription on these WAN-to-LAN links would exist. Therefore, these links should be protected with a queuing policy, as a queuing policy would be the only way to provide service level **guarantees** on these links, regardless of how rarely these queuing policies would engage. As discussed in the previous section, if LLQ is to be used for VoIP and TelePresence, then these should be configured with a dual-LLQ policy, similar to the simplified example provided in [Example 6-3](#) (but with different bandwidth values for both VoIP and TelePresence, based on how many calls of each type were being supported over these links).

Finally, WAN/VPN edge QoS policies would be required on all points labeled C in [Figure 6-5](#). The specifics of these WAN/VPN edge policy permutations are discussed in detail in this chapter.

TelePresence Branch LAN Edge

The LAN edge of the branch PIN performs essentially the same QoS services as does the campus access edge, namely the enforcement of a trust boundary, CoS-to-DSCP mapping (if required), optional TelePresence policing (to prevent network abuse of a trusted switch port), and queuing. However, there are some design considerations unique to the branch LAN edge discussed below.

TelePresence Branch LAN Edge QoS Design Considerations

Depending on the platform(s) used at the branch, QoS functions may be performed in hardware, in software, or in a combination of both. This is because Cisco IOS-based routers perform QoS in software, while Cisco Catalyst switches perform QoS in hardware. Additionally, some devices, such as the Cisco Integrated Services Routers, combine functionality from both product families within a single platform (for example, a Cisco ISR equipped with a Cisco EtherSwitch network module).

A general rule of thumb relating to QoS design is to **always enable QoS policies in hardware, rather than in software**, whenever a choice exists. This is because QoS policies performed in software require (marginal) incremental CPU loads to enforce (the actual incremental load varies according to platform, line rates, policy complexity, traffic patterns, and other variables). However, QoS policies performed in

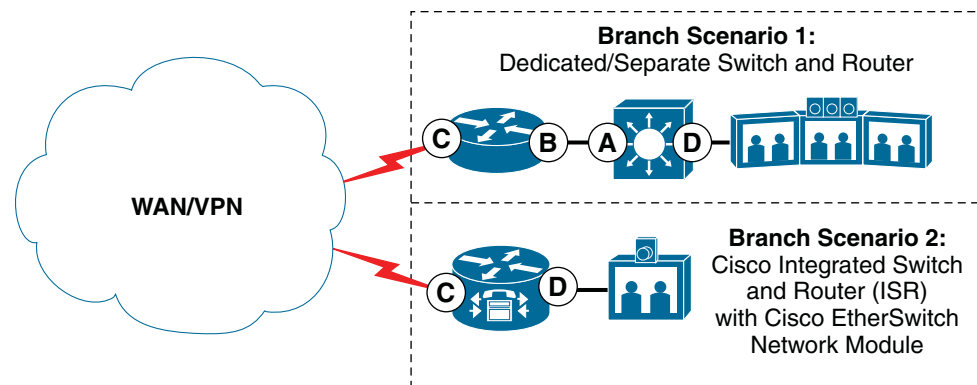
hardware are performed at line rates (GE or 10GE) without **any** incremental load to the CPU. Therefore, it is generally more efficient and effective to design QoS policies to be performed within hardware whenever possible.

Applying this rule of thumb to the branch LAN edge would typically result in one of two scenarios:

- The first scenario consists of a Cisco router on the WAN edge connecting to a Cisco Catalyst switch, which then connects to the branch endpoints, including the TelePresence system.
- The second scenario consists of a Cisco Integrated Services Router equipped with an EtherSwitch module performing all QoS functions within a single box.

These scenarios and their respective placements of QoS policies are illustrated in [Figure 6-6](#).

Figure 6-6 Enterprise Branch QoS Design Recommendations for TelePresence



Branch QoS Policies:

- | | |
|---|--|
| <p>(A) Interswitch Link Policies:
Trust DSCP
+ Queuing (CoS 4 and 5 → PQ)
+ Queuing (CoS 3 → Non-PQ)</p> | <p>(C) WAN/VPN Edge QoS Policies:
Trust for VoIP (EF)
+ LLQ or CBWFQ for TelePresence (CS4)
+ CBWFQ for Call-Signaling (CS3)</p> |
| <p>(B) Router LAN Edge Policies:
Trust DSCP (default)
+ (Optional) LLQ for VoIP (EF)
+ (Optional) LLQ for TelePresence (CS4)
+ (Optional) CBWFQ for Call-Signaling (CS3)</p> | <p>(D) Branch Access Edge Policies:
Trust for DSCP or Trust CoS
+ Map CoS 4 → DSCP CS4
+ Map CoS 5 → DSCP EF and CoS 3 → DSCP CS3
+ (Optional) Ingress Policing
+ Queuing (CoS 4 and 5 → PQ)
+ Queuing (CoS 3 → Non-PQ)</p> |

223247

We can see that for the most part the QoS policies deployed in the branch are reflective of the policies deployed at the campus WAN block. For example, the WAN/VPN edge QoS policies are applied to the branch router's WAN interface (to complement the WAN edge policies on the WAN/VPN aggregator's WAN edge). These are shown as points labeled C in [Figure 6-6](#). The considerations and details of these WAN/VPN edge policies are discussed in detail throughout this chapter.

In the case of a branch PIN using dedicated/separate switch(es) and router(s), the administrator has two additional policy points to configure: the router's LAN edge (shown the point labeled B in [Figure 6-6](#)) and the switch-to-router link (shown as the point labeled A in [Figure 6-6](#)). A queuing policy on the router's LAN edge (point B) is optional, as in many cases it may be theoretically impossible to congest this interface in the WAN-to-LAN direction: this interface would likely be a GigabitEthernet interface, and as such, well above the access rate of the WAN/VPN link. If it is to be configured with a queuing

policy, then a dual-LLQ policy would be recommended, such that VoIP and TelePresence are assigned to separate LLQs and Call-Signaling is protected with a CBWFQ. Such a policy would be similar to the simplified example provided in [Example 6-3](#) (but with different bandwidth values for both VoIP and TelePresence, based on how many calls of each type were being supported over these links).

Next, the administrator would need to configure the switch-to-router links, labeled as point A. These ports are essentially serving the same roles as campus interswitch links (as previously discussed in [Chapter 5, “Campus QoS Design for TelePresence”](#)). These ports should be configured to trust DSCP and perform hardware queuing, such that CoS 4 (TelePresence) and CoS 5 (VoIP) are assigned to the strict priority hardware queue and that CoS 3 (Call-Signaling) is assigned to a non-priority queue within the platform/linecard’s 1P×QyT queuing structure.

The final policy points are the branch access edges, which are shown as points labeled D in [Figure 6-6](#). These ports should be configured with either static or conditional trust of either DSCP or CoS (as described in detail in [Chapter 5, “Campus QoS Design for TelePresence”](#) in [Access Switch Port QoS Considerations](#)).

If CoS is trusted, then the necessary CoS-to-DSCP mappings must be in place, such that CoS 4 (TelePresence) is mapped to (the default value of) DSCP CS4 (32), CoS 5 (VoIP) is mapped to (the non-default value of) DSCP EF (46), and CoS 3 (Call-Signaling) is mapped to (the default value of) DSCP CS3 (24).

An optional recommendation for the branch access edge switch port connecting to a TelePresence primary codec is to configure a policer to prevent network abuse in case of a compromise of this trusted port. This recommendation helps prevent an unknowing and/or disgruntled individual that gains physical access to the TelePresence switch port from sending rogue traffic over the network that can hijack voice or video queues and easily ruin voice or video quality. Therefore, the administrator may choose to limit the scope of damage that such network abuse may present by configuring access edge policers on TelePresence switch ports to drop (or remark to Scavenger - DSCP CS1) out-of-profile traffic originating on these ports. This is not only a Cisco recommended best practice, but is also reflected in RFC 4594, which recommends edge policing the Real-Time Interactive service class via a single-rate policer. If such a policer is configured, it is recommended to use Per-Port/Per-VLAN policers, whenever supported. In this manner, a set of policers may be applied to the Voice VLAN to ensure that voice, video, and call signaling traffic are performing within normal levels and a separate, more stringent policer can be applied to the data VLAN.


Note

Access Edge policers are described in detail on a platform-by-platform basis in Chapter 2, “Campus QoS Design” of the QoS SRND at www.cisco.com/go/srnd.

Additionally, the recommended burst parameter for TelePresence policers is discussed in detail in [TelePresence Branch WAN Edge LLQ Policy](#).

Finally, it is recommended to enable queuing on these branch access edge links in the case of congestion. While the likelihood of such an event is rare, these may occur during DoS/worm attacks; therefore, provisioning queuing policies on these links is mandatory to provide service level guarantees in **any** event.

In the case of a branch PIN using an ISR with an EtherSwitch module, the only real twist is that the administrator would configure the WAN/VPN edge policies in the router console mode, and then switch to the EtherSwitch console mode (using the **service-module gigabitEthernet module/number session** IOS command) to configure hardware QoS policies on the EtherSwitch module (which is essentially a Cisco Catalyst 3750 switch), using the Catalyst 3750 configuration recommendations provided in [Chapter 5, “Campus QoS Design for TelePresence.”](#)

TelePresence Branch LAN Edge QoS Designs

The configuration details for branch LAN edge policies are identical to the platform- and linecard-specific policies used in the campus access edge, which have already been covered in detail in [Chapter 5, “Campus QoS Design for TelePresence”](#) and as such it would be redundant to again detail these designs in this chapter.

TelePresence Branch WAN Edge

The WAN edge of the branch is likely the most congested node within the network and as such requires the most attention from a QoS perspective. Let us now recap some of the considerations of the WAN edge and then delve into design detail.

TelePresence Branch WAN Edge Design Considerations

In a private WAN design, the principal decision that the administrator needs to make has already been covered in detail, namely whether to service TelePresence with a LLQ or a CBWFQ. A point to keep in mind with respect to private WAN scenarios is that additional costs are typically not incurred when provisioning additional traffic with strict priority servicing and, as such, these scenarios are generally more conducive to provisioning TelePresence in a LLQ than other VPN scenarios.

Once the LLQ versus CBWFQ decision is made, then the WAN edge policies are fairly straightforward and are a function of this decision coupled with the number of traffic classes that have been defined by the enterprise’s strategic business QoS objectives (as discussed in [Chapter 4, “Quality of Service Design for TelePresence”](#)).

TelePresence Branch WAN Edge QoS Design

To recap, LLQ provides superior levels of service for TelePresence as compared to CBWFQ, yet may entail additional costs or other constraints. Therefore, each administrator must make an informed decision as to which WAN edge queuing strategy (LLQ or CBWFQ) to employ. Both configuration options are presented in detail in the following sections.

TelePresence Branch WAN Edge LLQ Policy

If TelePresence is to be assigned to an LLQ, then in addition to adequately provisioning priority bandwidth to the LLQ, one additional design parameter needs to be calculated: the burst parameter of the implicit policer of the LLQ.

The implicit policer for the LLQ is a token-bucket algorithm policer (like any other IOS or Catalyst policer) and as such needs a burst parameter to be defined in order to police to a sub-line rate. To better understand why this is so, let us briefly recap how token-bucket policers work.

To regulate transmissions at sub-line rates, the concept of an interval must be applied. An interval is a sub-second period of time during which an application may send traffic. For example, if a policer was to limit an application to 15 Mbps of a 45 Mbps circuit, then the policer would allow the application to transmit for a total of 333 ms per second (15 / 45 Mbps) and it would drop any packets offered during the remaining 667 ms per second. Now, if the application sent all its traffic in a single burst, this could tie up the circuit for up to one-third of a second, which may cause excessive jitter and/or drops to other applications. Therefore, rather than allowing a single interval per second for an application to send

traffic, it is generally more efficient to configure policers that allow for transmission over multiple sub-second intervals. The amount of traffic that an application can transmit during a sub-second interval is called the committed burst or Bc. The time interval itself is referred to as the time constant or Tc. The relationship between the burst, the interval, and the overall policing rate is:

$$\text{Bc} = \text{Policing Rate} * \text{Tc}$$

Now let us apply this theory to TelePresence traffic patterns so that we can define an optimal value for the burst parameter of the LLQ's implicit policer.

TelePresence codecs, whether operating at 720p or 1080p resolution, display 30 frames per second. Put differently, TelePresence codecs send information representing one frame every 33 ms (1 second/30 frames-per-second). We can use this information as a starting point, as it directly correlates to our interval (Tc).

Now, if TelePresence had a fixed packet size and a constant packetization rate (like VoIP), then we could simply divide the per second bandwidth requirements (shown in Table 4-1 in Chapter 4, “Quality of Service Design for TelePresence”) by 30 (fps) to arrive at our burst parameter. For example, under this assumption, a CTS-3000 transmitting at 1080p-Best would need a burst parameter of 62,500 Bytes (15 Mbps / 30 fps / 8 bits).

**Note**

However, while a fixed packet size and a constant packetization rate might make burst calculations a bit simpler, these would result in exponentially higher bandwidth requirements for TelePresence. As discussed in Chapter 4, “Quality of Service Design for TelePresence,” if TelePresence were uncompressed, it would result in 1.5 Gbps of bandwidth per display, rendering TelePresence virtually undeployable—especially over wide area networks.

However, TelePresence does not have a fixed packet size, nor a constant packetization rate, as it utilizes advanced video compression techniques to achieve compression rates of over 99%, thus massively reducing the bandwidth requirements for TelePresence and rendering it more deployable, even over WANs. Notwithstanding, these high compression algorithms within TelePresence systems do have a direct impact in burst calculations for policers, such as the implicit policer within LLQ.

Therefore, to configure the policing burst such that it does not drop TelePresence traffic, we have to analyze what would be the maximum transmission (in Bytes) within a 33 ms interval—in other words, the worst-case scenario per frame of TelePresence video. In H.264 video, which TelePresence systems utilize, this worst-case scenario would be the full screen of (spatially-compressed) video, which is periodically sent, known as the Instantaneous Decoding Refresh (IDR) frame. The IDR frame is the key frame that subsequent video frames reference, sending only differential information between subsequent frames and the IDR frame, rather than the full-picture again.

**Note**

For more information about H.264 video encoding, refer to RFC 3964 “RTP Payload Format for H.264 Video” at <http://www.ietf.org/rfc/rfc3984>.

The maximum IDR frame sizes observed during extensive testing of TelePresence systems (using CTS software version 1.1.0 [256D]) was 64 KB. Therefore, the LLQ burst parameter should be configured to permit up to 64 KB of burst per frame per screen. In the case of a triple-display CTS-3000 systems, we should allow for 192 KB of burst (3 * 64 KB) in the rare event of a “triple-IDR storm,” where all three codecs send IDR frames simultaneously.

**Note**

If Cisco design recommendations for TelePresence room lighting and other environmental variables are not followed, then IDR frame sizes may vary in size beyond 64 KB, which may in turn affect the network QoS policies.

However, it bears mentioning that the version of CTS software used in this phase of testing (1.1.0 [256D]) did not support an auxiliary video stream. If newer versions of CTS software are being used or if the use of an auxiliary video stream (for sharing PowerPoint presentations, etc.) is planned, then a larger value of TelePresence burst would be required. Subsequent testing has shown that a value of 256 KB is sufficient to support TelePresence with an auxiliary video stream (192 KB for worst-case primary video + 64 KB for worst-case auxiliary video). Therefore, the examples that follow utilize this higher burst value to adequately provision for the use of TelePresence with an auxiliary video stream; if, on the other hand, such use is not planned, then a value of 192 KB is sufficient for TelePresence burst provisioning.

Now let us put this all together into a configuration. To quickly recap, the full syntax of the LLQ command in Cisco IOS is:

```
priority { bandwidth-kbps | percent percentage } [burst]
```

As can be seen, the burst parameter is an optional parameter that can be explicitly defined as part of the **priority** command. If the burst is not explicitly defined, then it defaults to a value computed as 200 ms of traffic at the configured LLQ bandwidth rate. However, it is important to note that the burst value is expressed in Bytes (not bits).

For example, if **priority 1000** was configured for a class, then the default burst parameter would be set to 25000 Bytes (1000 kbps * 200 ms / 8 bits). This value would not appear in the configuration, but could be verified with a **show policy map interface** verification command.

Let us look at the worst-case burst for a CTS-1000 system. Applying the IDR as the worst-case burst scenario for TelePresence primary video (at 64 KB) coupled with an allowance of auxiliary video bursting of the same amount (64 KB), the configuration to provision a branch WAN edge queuing policy that provisions a TelePresence CTS-1000 system running at 1080p-Best (with the optional support of an auxiliary video stream) to a LLQ with an optimal burst parameter (of 128 KB) is shown in [Example 6-4](#).

Example 6-4 Dual-LLQ Branch WAN Edge Policy for VoIP and TelePresence (CTS-1000 at 1080p-Best with Auxiliary Video)

```
policy-map WAN-EDGE
  class VOIP
    priority percent 10           ! LLQ for VoIP (example amount of BW)
  class TELEPRESENCE
    priority 5500 128000        ! LLQ for CTS-1000 (1080p-Best + aux video)
  class DATA
    ...
```

Likewise, the configuration to provision a branch WAN Edge queuing policy that provisions a TelePresence CTS-3000 system running at 1080p-Best (with the optional support of an auxiliary video stream) to a LLQ with an optimal burst parameter (of 256 KB) is shown in [Example 6-5](#).

Example 6-5 Dual-LLQ Branch WAN Edge Policy for VoIP and TelePresence (CTS-3000 at 1080p-Best with Auxiliary Video)

```
policy-map WAN-EDGE
  class VOIP
    priority percent 10           ! LLQ for VoIP (example amount of BW)
  class TELEPRESENCE
    priority 15000 256000       ! LLQ for CTS-3000 (1080p-Best + aux video)
  class DATA
    ...
```

These configurations can be verified with the following command:

- **show policy-map interface**

TelePresence Branch WAN Edge CBWFQ Policy

If, on the other hand, TelePresence is to be assigned to a CBWFQ, then in addition to adequately provisioning guaranteed bandwidth to the CBWFQ, one additional design parameter needs to be considered, the length of the CBWFQ.

By default, Class-Based Weighted Fair Queues are 64 packets deep. Extensive testing has shown that this default queue-depth has at times resulted in tail-drops when provisioned to protect TelePresence flows. Therefore, on most interfaces it is recommended to increase the default queue-depth for the TelePresence queue to 128 packets, using the **queue-limit 128** command in conjunction with the CBWFQ **bandwidth** command.

An example policy provisioning a TelePresence CTS-3000 system running at 1080p-Best (with the optional support of an auxiliary video stream) to a CBWFQ, with an extended queue-depth to 128 packets, is shown in [Example 6-6](#).

Example 6-6 CBWFQ Branch WAN Edge Policy for TelePresence (CTS-3000 at 1080p-Best with Auxiliary Video)

```
policy-map WAN-EDGE
  class VOIP
    priority percent 10           ! LLQ for VoIP (example amount of BW)
  class TELEPRESENCE
    bandwidth 15000              ! CBWFQ for CTS-3000 (1080p-Best + aux video)
    queue-limit 128              ! Extended queue-limit for TelePresence CBWFQ
  class DATA
    ...
```

This configuration can be verified with the following command:

- **show policy-map interface**

TelePresence Branch T3/DS3 WAN Edge Design

When configuring a WAN edge policy for TelePresence, there are a couple of additional considerations that need to be taken into account when using T3 interfaces, namely, adjusting the hold-queue size (if needed) to accommodate all LLQ/CBWFQs and tuning the Tx-Ring to minimize TelePresence jitter on converged links.

The total number of buffers that the IOS software allocates for queuing per interface (regardless of whether the interface is configured with FIFO, LLQ, or CBWFQ) is called the output queue or hold-queue. The size of the output queue and can be adjusted with the **hold-queue** interface command to a value between 0 and 4096 packets; the default output queue size of a T3 serial interface is 1000 packets.

Normally the default hold-queue size is sufficient for a T3 interface. Consider our worst-case example, where we have a 12-class RFC 4594-based QoS policy and let us choose the option with TelePresence assigned to a CBWFQ. Besides TelePresence, there are 10 CBWFQs, each a default queue-depth of 64 packets, for an output queue depth of, so far, 640 packets. Additionally, there is the 128 packet extended queue-depth for the TelePresence CBWFQ, bringing our running total output queue depth to 768 packets which, even factoring a moderate allowance for LLQ queue depth, is well below our default value of 1000 packets for this T3 interface.

However, should the administrator—for whatever reason—require expanding the queue depths of each CBWFQ to 128, then the output queue depth requirement would be at least (11 * 128) 1408 packets, not even factoring a moderate allowance for the LLQ. In such a scenario, LLQ traffic could be impacted if

the output queue size was not expanded accordingly. For instance, in such a case, the network administrator could increase the size of the output queue to 1500 packets by using the **hold-queue 1500** interface command.

Earlier in this chapter we introduced the Tx-Ring. To quickly recap, the Tx-Ring represents the size of the final output buffer (a FIFO queue) that maximizes physical link bandwidth utilization by matching the outbound packet rate on the router with the physical interface rate. The Tx-Ring also serves to indicate interface congestion to the IOS software. Prior to interface congestion, packets are sent on a FIFO basis to the interface via the Tx-Ring. However, when the Tx-Ring fills to its queue-depth/limit, then it signals to the IOS software to engage any LLQ/CBWFQ policies that have been attached to the interface. Subsequent packets are then queued within IOS according to these LLQ/CBWFQ policies, dequeued into the Tx-Ring, and then sent out the interface in a FIFO manner. These operations are illustrated in [Figure 6-2](#), [Figure 6-3](#), and [Figure 6-4](#).

The Tx-Ring can be configured on certain platforms, such as the Cisco PA-T3+ port adapter interface, with the **tx-ring-limit** interface command. The value of the **tx-ring-limit** number can be from 1 to 32,767 packets. The default for serial interfaces on the PA-T3+ is 64 packets.

During testing it was observed that the default tx-ring-limit limit of 64 packets was shown to cause somewhat higher jitter values to TelePresence traffic during fully-congested scenarios. The reason for this is the bursty nature of TelePresence traffic. Even though TelePresence traffic is prioritized when LLQ/CBWFQ policies are active, if there are no TelePresence packets to send, the FIFO Tx-Ring is filled with other traffic. When a new TelePresence packet arrives, even if it gets priority treatment from the Layer 3 LLQ/CBWFQ queuing system, the packet are dequeued into the FIFO Tx-Ring when space is available. However, with the default settings, there can be as many as 63 (non-TelePresence) packets in the Tx-Ring in front of that TelePresence packet. In such a worst-case scenario it would take as long as 17 ms to transmit these non-TelePresence packets out of the T3 interface. This 17 ms of instantaneous delay (i.e., jitter) exceeds the jitter target for TelePresence.

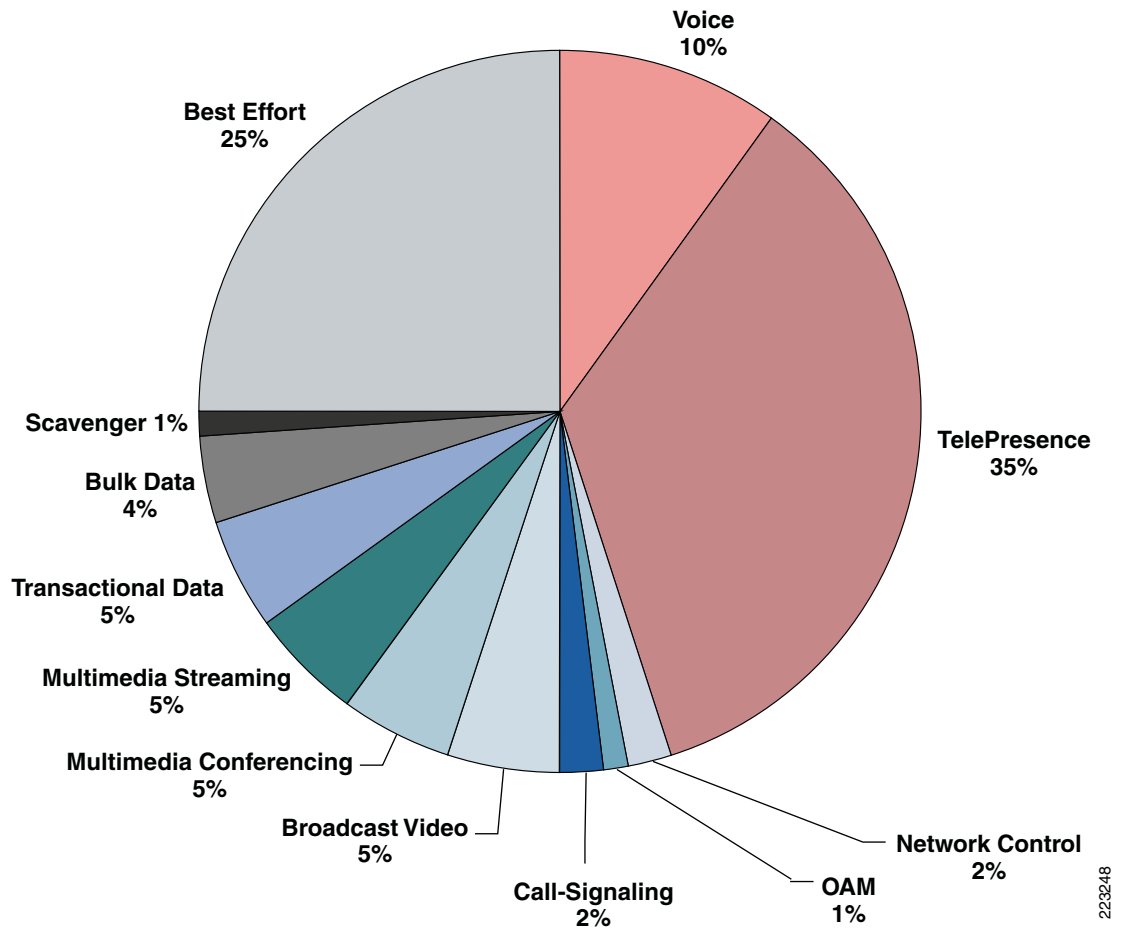
Additionally, a lower Tx-Ring results in the IOS software engaging congestion management policies sooner and more often, resulting in lower overall jitter values for priority traffic, such as TelePresence. On the other hand, setting the value of the Tx-Ring too low may result in significantly higher CPU utilization rates, as the processor is being continually interrupted to engage queuing policies, even when congestion rates are just momentary bursts and not sustained rates. Thus when tuning the Tx-Ring, a trade-off setting is required such that jitter is minimized, but not at the expense of excessive CPU utilization rates. Therefore, extensive testing has shown that setting the Tx-Ring to a value of 10 packets is optimal on converged T3 links supporting TelePresence and other applications (such as voice, data, and other video applications). This can be achieved by using the **tx-ring-limit 10** interface command.

**Note**

As explained above, tuning the Tx-Ring to value of 10 is only required on converged links that support additional applications beyond TelePresence. On T3 circuits that are dedicated to TelePresence, lowering the Tx-Ring to a non-default value is not required; in fact, such tuning can actually deteriorate the quality of TelePresence calls on such dedicated T3 circuits. It is important to keep in mind that in the case of dedicated circuits, it is not other applications that could potentially fill the Tx-Ring, but rather other TelePresence flows. Furthermore, when properly provisioned, dedicated links should not generate sustained congestion scenarios.

Now let us put this all together into a full example. In this 12-class RFC 4594-based case-study example, it was decided to service TelePresence traffic over the branch T3 WAN edge in a CBWFQ, as provisioning both VoIP and TelePresence in a dual-LLQ design would, in this case, require 45% of priority queuing, which would cause excessive variations in application response times to the other 10 application classes. The WAN edge bandwidth allocations for this case-study example are shown in [Figure 6-7](#).

Figure 6-7 Case Study Example Bandwidth Allocations of a RFC 4594-Based LLQ/CBWFQ Policy Over a Branch T3 WAN Edge



The corresponding configuration for this case study example is shown in [Example 6-7](#).

Example 6-7 Case Study Example Configuration of a RFC 4594-Based LLQ/CBWFQ Policy (with TelePresence in a CBWFQ) Over a Branch T3 WAN Edge

```

!
class-map match-all VOICE
  match dscp ef
class-map match-all TELEPRESENCE
  match dscp cs4
class-map match-all NETWORK-CONTROL
  match dscp cs6
class-map match-all OAM
  match dscp cs2
class-map match-all CALL-SIGNALING
  match dscp cs3
class-map match-all BROADCAST-VIDEO
  match dscp cs5
class-map match-all MULTIMEDIA-CONFERENCING
  match dscp af41 af42 af43
class-map match-all MULTIMEDIA-STREAMING
  match dscp af31 af32 af33
class-map match-all TRANSACTIONAL-DATA
  match dscp af21 af22 af23
  
```

223248

```

class-map match-all BULK-DATA
  match dscp af11 af12 af13                                ! Bulk-Data markings
class-map match-all SCAVENGER
  match dscp cs1                                          ! Scavenger marking
!
!
policy-map WAN-EDGE-T3
  class VOICE
    priority percent 10                                  ! LLQ for VoIP
  class TELEPRESENCE
    bandwidth percent 35                                ! CBWFQ for TP (CTS-3000)
    queue-limit 128                                    ! Expanded Queue-Limit for TP
  class NETWORK-CONTROL
    bandwidth percent 2                                  ! CBWFQ for Routing
  class OAM
    bandwidth percent 1                                  ! CBWFQ for Ops/Admin/Mgmt
  class CALL-SIGNALING
    bandwidth percent 2                                  ! CBWFQ for Call-Signaling
  class BROADCAST-VIDEO
    bandwidth percent 5                                  ! CBWFQ for Broadcast Video
  class MULTIMEDIA-CONFERENCING
    bandwidth percent 5                                  ! CBWFQ for IP/VC
    random-detect dscp-based                            ! DSCP-WRED for IP/VC
  class MULTIMEDIA-STREAMING
    bandwidth percent 5                                  ! CBWFQ for Streaming-Video
    random-detect dscp-based                            ! DSCP-WRED for Stream-Video
  class TRANSACTIONAL-DATA
    bandwidth percent 5                                  ! CBWFQ for Trans-Data
    random-detect dscp-based                            ! DSCP-WRED for Trans-Data
  class BULK-DATA
    bandwidth percent 4                                  ! CBWFQ for Bulk Data
    random-detect dscp-based                            ! DSCP-WRED for Bulk Data
  class SCAVENGER
    bandwidth percent 1                                  ! Minimum CBWFQ for Scavenger
  class class-default
    bandwidth percent 25                                  ! CBWFQ for Best Effort
    random-detect                                        ! WRED for Best Effort
!
...
!
interface Serial6/0
  description BRANCH-TO-CAMPUS-T3
  ip address 192.168.2.9 255.255.255.252
  tx-ring-limit 10                                     ! Tuned T3 Tx-Ring
  dsu bandwidth 44210
  framing c-bit
  cablelength 10
  serial restart-delay 0
  max-reserved-bandwidth 100                           ! LLQ/CBWFQ BW Override
  service-policy output WAN-EDGE-T3                   ! Attaches policy to T3 int
!

```

**Note**

Since this policy supports a less-than Best Effort Scavenger-class, it requires an explicit CBWFQ to be configured on class-default (the Best Effort class); otherwise, the queuing algorithm robs bandwidth from class-default to service Scavenger traffic. Additionally, when class-default is configured with an explicit **bandwidth** command, then the **max-reserved-bandwidth** interface command must also be configured on the outgoing interface. Additional details on this behavior can be found in the QoS SRND at www.cisco.com/go/srnd on pages 3-8 and 3-9.

Optionally, if the network administrator chooses to use LLQ instead of CBWFQ, then the only change to the above policy would be to the TELEPRESENCE class within the WAN-EDGE-T3 policy-map, as shown in [Example 6-8](#).

Example 6-8 Policy Amendment to Example 6-7 to Provision TelePresence in a Dual-LLQ Policy Over a T3 Branch WAN Edge

```
!
policy-map WAN-EDGE-T3
  class VOICE
    priority percent 10           ! LLQ for VoIP
  class TELEPRESENCE
    priority percent 35 256000   ! LLQ for CTS-3000 (1080p-Best + aux video)
  class NETWORK-CONTROL
...

```

These configurations can be verified with the following commands:

- **show interface**
- **show policy-map interface**
- **show controllers Serial *module/interface* | include tx_limited**

TelePresence Branch OC3-POS WAN Edge Design

When configuring a WAN edge policy for TelePresence, there are a couple of additional considerations that need to be taken into account when using OC3-POS interfaces, such as the Cisco 7600 SPA-2XOC3-POS. Both of these considerations relate to Low-Latency Queuing:

- There is a 35% hard limit to the amount of traffic that can be configured with priority queuing.
- There is different configuration syntax for enabling LLQ on these interfaces.

Let us discuss these considerations in more detail.

The first consideration is fairly straightforward: on these OC3-POS interfaces, there is a hard limit of 35% for the sum of all traffic that can be configured with LLQ. This means that either a single LLQ class can be configured with a maximum of 54.25 Mbps or the sum of all LLQ classes can be configured for a combined maximum of 54.25 Mbps. This hard-limit, incidentally, is quite consistent with Cisco's "33% LLQ Rule."

The second consideration has to do with a change in syntax for configuration of LLQ on these OC3-POS interfaces. The configuration syntax for strict priority queuing on an OC3 POS interface is such the **priority** command does not include a bandwidth parameter, either as an absolute value (defined in kbps) or as a percentage of the link's bandwidth. That being said, there is correspondingly no implicit policer within the LLQ **priority** command, but rather an explicit policer must be configured on the policy-map class and then the **priority** command can be applied to the class. This difference is not limited to syntax only, but also affects behavior, as an implicit policer only engages when the LLQ is active (i.e., during periods of congestion), but an explicit policer, such as required for a LLQ class on an OC3-POS interface, is always on.

**Note**

The always-on nature of explicit policers is advantageous from an Admission Control perspective. For example, without a comprehensive, network-aware Call Admission Control system in place, there would be no way to always enforce limits on TelePresence traffic without explicit policers (remember, implicit policers, like those included within LLQ, are only active during congestion scenarios). Admission Control considerations and designs are discussed in more detail in [Chapter 8, “Capacity Planning and Call Admission Control.”](#)

The configured explicit LLQ policer may be either a single-rate or a dual-rate policer. When applied to TelePresence traffic, only a single-rate policer is relevant (as we are not interested in marking down excess TelePresence traffic, which we could do with two levels of granularity via a dual-rate policer). Additionally, testing has shown that using a dual-rate policer offers no performance advantages whatsoever; therefore it is recommended to configure a single-rate policer on the TelePresence LLQ class.

As with an implicit policer, a committed burst parameter is required when defining an explicit policer. As discussed in [TelePresence Branch WAN Edge LLQ Policy](#), the recommended value for TelePresence committed burst for a CTS-3000 system running 1080p-Best (with optional auxiliary video support) is 256 KB.

Reflecting the foregoing points, the configuration syntax for creating a single-rate explicit policer and applying it to a TelePresence LLQ class (for a CTS-3000 system running at 1080p-Best with auxiliary video) LLQ is shown in [Example 6-9](#).

Example 6-9 TelePresence LLQ Policy Over a OC3-POS Branch WAN Edge

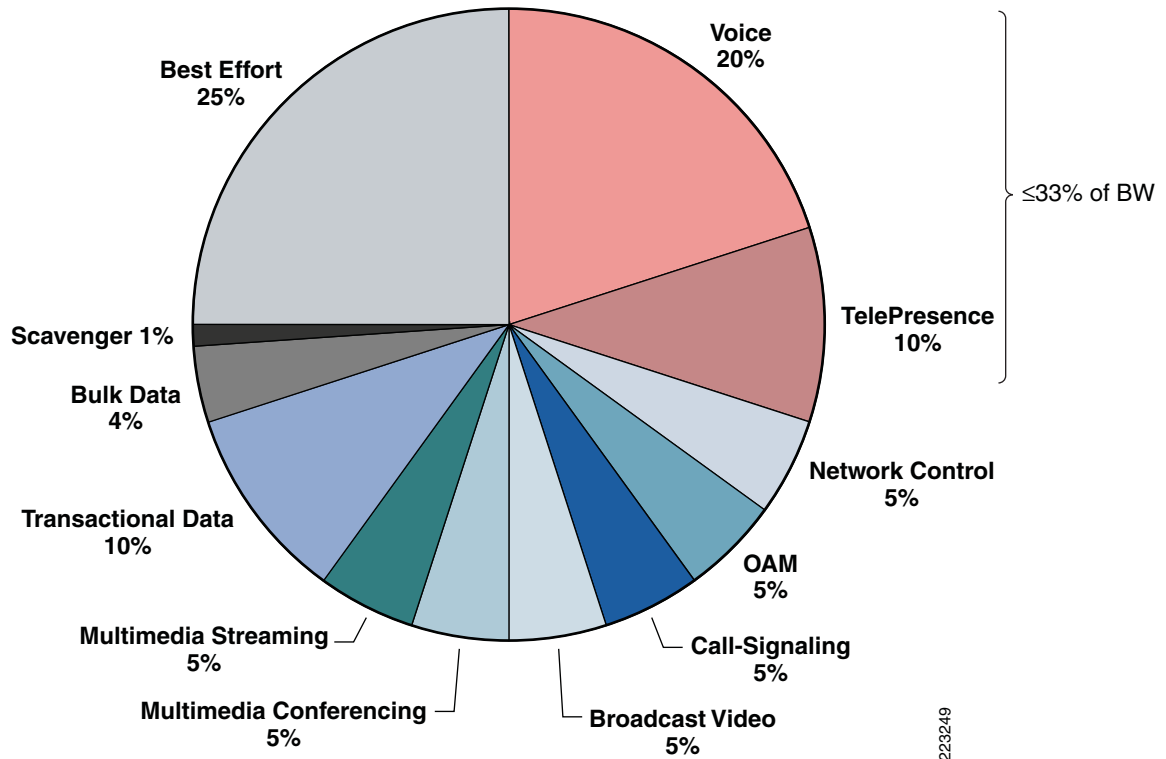
```
!
policy-map WAN-EDGE
class TELEPRESENCE
  police cir 15000000           ! TP is policed to 15 Mbps
    bc 256000                 ! Bc is 256 KB
    conform-action transmit    ! Conforming action --> transmit
    exceed-action drop        ! Single-Rate Policing action
  priority                    ! LLQ command for OC3-POS
!
```

This configuration can be verified with the following command:

- **show policy-map interface**

Now let us again put this all together into a full example. In this 12-class RFC 4594-based case-study example, it has been decided to service TelePresence traffic over the branch OC3-POS WAN edge in a dual-LLQ design, along with voice. The WAN edge bandwidth allocations for this case-study example are shown in [Figure 6-8](#).

Figure 6-8 Case Study Example Bandwidth Allocations of a RFC 4594-Based LLQ/CBWFQ Policy Over a Branch OC3-POS WAN Edge



The corresponding configuration for this second case study example is shown in [Example 6-10](#) (the class-maps are not repeated, as these do not change).

Example 6-10 Case Study Example Configuration of a RFC 4594-Based LLQ/CBWFQ Policy (with TelePresence in a Dual-LLQ) Over a Branch OC3-POS WAN Edge

```

!
policy-map WAN-EDGE-OC3-POS
  class VOICE
    police cir 31000000           ! Voice is policed to 31 Mbps (20%)
      bc 15500                   ! Bc is 15.5 KB
      conform-action transmit     ! Conforming action --> transmit
      exceed-action drop         ! Single-Rate Policing action
      priority                   ! LLQ command for OC3-POS
  class TELEPRESENCE
    police cir 15000000         ! TP is policed to 15 Mbps
      bc 256000                 ! Bc is 256 KB
      conform-action transmit     ! Conforming action --> transmit
      exceed-action drop         ! Single-Rate Policing action
      priority                   ! LLQ command for OC3-POS
  class NETWORK-CONTROL
    bandwidth percent 5        ! CBWFQ for Routing
  class OAM
    bandwidth percent 5        ! CBWFQ for Network Management
  class CALL-SIGNALING
    bandwidth percent 5        ! CBWFQ for Call-Signaling
  class BROADCAST-VIDEO
    bandwidth percent 5        ! CBWFQ for Broadcast Video
  class MULTIMEDIA-CONFERENCEING
    bandwidth percent 5        ! CBWFQ Video-Conferencing

```



```

    random-detect dscp-based                ! DSCP-WRED for Video-Conferencing
class MULTIMEDIA-STREAMING
    bandwidth percent 5                     ! CBWFQ for Streaming-Video
    random-detect dscp-based                ! DSCP-WRED for Streaming-Video
class TRANSACTIONAL-DATA
    bandwidth percent 10                   ! CBWFQ for Transactional Data
    random-detect dscp-based                ! DSCP-WRED for Transactional Data
class BULK-DATA
    bandwidth percent 4                     ! CBWFQ for Bulk Data
    random-detect dscp-based                ! DSCP-WRED for Bulk Data
class SCAVENGER
    bandwidth percent 1                     ! Minimum CBWFQ for Scavenger
class class-default
    bandwidth percent 25                   ! CBWFQ for Best Effort
    random-detect! WRED for Best Effort
!
...
interface POS3/0/1
description BRANCH-TO-CAMPUS-OC3-POS
ip address 192.168.5.1 255.255.255.252
clock source internal
service-policy output WAN-EDGE-OC3-POS ! Attaches policy to OC3-POS
!

```

**Note**

No Tx-Ring tuning is required on the OC3-POS link; neither is a **max-reserved-bandwidth 100** interface command required.

Optionally, if the network administrator chooses to use CBWFQ instead of LLQ for TelePresence, then the only change to the above policy would be to the TELEPRESENCE class within the WAN-EDGE-OC3-POS policy-map, as shown in [Example 6-11](#).

Example 6-11 Policy Amendment to Example 6-10 to Provision TelePresence in a CBWFQ Over an OC3-POS Branch WAN Edge

```

!
policy-map WAN-EDGE-OC3-POS
class VOICE
    police cir 31000000                     ! Voice is policed to 31 Mbps (20%)
    bc 15500                                ! Bc is 15.5 KB
    conform-action transmit                 ! Conforming action --> transmit
    exceed-action drop                      ! Single-Rate Policing action
    priority                                ! LLQ command for OC3-POS
class TELEPRESENCE
    bandwidth percent 10                   ! CBWFQ for TelePresence
class NETWORK-CONTROL
...

```

**Note**

Testing has shown that extending the TelePresence CBWFQ queue-limit beyond the default value of 64 packets is not required on OC3-POS interfaces because of the extremely fast serialization rate—relative to TelePresence transmission rates—of these interfaces.

These configurations can be verified with the following command:

- **show policy-map interface**

TelePresence Branch IPSec VPN Edge

In the initial releases of Cisco TelePresence software, native encryption within the codecs was not supported. Nonetheless, for certain enterprises, encryption is a business requirement for all IP communications. This business requirement can be achieved by performing IPSec encryption and decryption within the network infrastructure, such as at the branch WAN/VPN edges over a private WAN or an MPLS VPN infrastructure. However, it is good to review some design considerations relating to IPSec and QoS interaction prior to examining the configuration details required for these deployments.



Note

Cisco does not recommend deploying TelePresence over the Internet—with or without IPSec encryption—as critical service level parameters, such as latency, jitter, and loss, cannot be guaranteed over the Internet. These IPSec designs for TelePresence are intended for use over private WAN and/or MPLS VPN scenarios.

TelePresence Branch IPSec VPN Edge Considerations

One of the first considerations is that IPSec adds network overhead to the packets. How much overhead depends on the encryption and tunneling options defined within the security associations. For example, a typical IPSec configuration uses IPSec Tunnel Mode with **esp-3des** and **esp-md5-hmac**, which results in an overhead of 56 bytes, accounted for as shown in [Table 6-1](#).

Table 6-1 *IPSec Network Overhead Breakdown*

Component	Overhead in Bytes
IPSec header (bytes)	20
ESP Header (bytes SPI)	4
ESP Header (bytes Sequence)	4
IOS ESP-DES/3DES (bytes IV)	8
ESP-DES/3DES 64-bit (bytes pad)	6
ESP Trailer (byte PAD length)	1
ESP Trailer (byte Next Header)	1
ESP MD5 96 digest (bytes)	12
Total IPSec Overhead	56

This being the case, since the average packet size for TelePresence is around 1200 Bytes, encryption overhead is typically <5%. [Table 6-2](#) shows a detailed breakdown of the respective encrypted bandwidth requirements for all TelePresence motion-handling and resolution options.

Table 6-2 TelePresence Bandwidth Requirements with IPSec Encryption

Motion Handling	Best	Better	Good	Best	Better	Good
Resolution	1080p	1080p	1080p	720p	720p	720p
CTS 1000						
Max with IPSec overhead (Kbps)	5,792	5,194	4,596	4,596	3,400	2,204
CTS 3000						
Max with IPSec overhead (Kbps)	15,360	13,566	11,772	11,772	8,184	4,596

Another important consideration is the interaction of IPSec and QoS, particularly with respect to Anti-Replay. In order to understand this interaction implication, it is beneficial to briefly recap the purpose and function of IPSec Anti-Replay.

IPSec offers inherent message-integrity mechanisms to provide a means to identify whether an individual packet is being replayed by an interceptor or hacker. This concept is called connectionless integrity. IPSec also provides for partial sequence integrity, preventing the arrival of duplicate packets.

**Note**

Anti-Replay concepts are outlined in RFC 2401, “Security Architecture for the Internet Protocol” at www.ietf.org/rfc/rfc2401.

When ESP authentication is configured in an IPSec transform set, for each security association, the receiving IPSec peer verifies that packets are received only once. Because two IPSec peers can send millions of packets, a 64-packet sliding window is implemented to bind the amount of memory required to tally the receipt of a peer’s packets. Packets can arrive out of order, but they must be received within the scope of the window to be accepted. If they arrive too late (outside the window), they are dropped.

The operation of the Anti-Replay window protocol is as follows:

1. The sender assigns a unique sequence number (per security association) to encrypted packets.
2. The receiver maintains a 64-packet sliding window, the right edge of which includes the highest sequence number received.
3. The receiver evaluates the received packet’s sequence number:
 - If a received packet’s sequence number falls within the window and was not received previously, the packet is accepted and marked as received.
 - If the received packet’s sequence number falls within the window and previously was received, the packet is dropped and the replay error counter is incremented.
 - If the received packet’s sequence number is greater than the highest sequence in the window, the packet is accepted and marked as received, and the sliding window is moved “to the right.”
 - If the received packet’s sequence number is less than the lowest sequence in the window, the packet is dropped and the replay error counter is incremented.

While Anti-Replay is useful in validating message integrity, in a converged IPSec VPN implementation with QoS enabled, lower-priority packets are often delayed so that higher-priority packets receive preferential treatment, which has the unfortunate side effect of sufficiently reordering packets so they are out of sequence from an IPSec Anti-Replay perspective. Therefore, there is a concern that through the normal QoS prioritization process, the receiver might drop packets as Anti-Replay errors, when, in fact, they are legitimately sent or received packets.

Traffic assigned to CBWFQ classes is much more sensitive to Anti-Replay than traffic assigned to a LLQ. This is because LLQ traffic is always sent in order, with strict priority; but CBWFQ traffic may be delayed by other CBWFQ flows and be sent in gaps exceeding the receiver's 64-packet sliding Anti-Replay window. Furthermore, by default, each CBWFQ class receives a queue with a length of 64 packets. Meanwhile, the receiving IPSec peer has a single 64-packet Anti-Replay window (per IPSec Security Association) with which to process packets from **all** LLQ and CBWFQ bandwidth classes. Therefore, a mismatch is created between the queue depths on the sender's output interface (multiple queues of 64 packets each) as compared to the width of the receiver's Anti-Replay window (a single sliding window of 64 packets per SA). As more bandwidth classes are defined in the sender's policy map, this mismatch increases. This is an inefficient use of expensive WAN/VPN bandwidth, as many packets are transmitted only to be dropped before decryption.

Cisco IOS allows the Anti-Replay window to be expanded (up to a maximum value of 1024 packets) or, alternatively, to be disabled entirely.

During testing it was observed that when TelePresence was provisioned within a dual-LLQ design, with a default-sized Anti-Replay window (of 64 packets), Anti-Replay errors did not affect either TelePresence or voice flows; however, there were significant Anti-Replay errors occurring on CBWFQ classes, inline with the behavior described above. These errors were reduced as the Anti-Replay window was enlarged, to the maximum of 1024 packets, yet were only eliminated altogether when Anti-Replay was disabled.

Additionally, when TelePresence was provisioned with a CBWFQ, with a default-sized Anti-Replay window, significant replay errors occurred on TelePresence flows, resulting in unusable call-quality. Replay errors were still noticed even when the Anti-Replay sliding window was set to the maximum of 1024 packets and were only eliminated when the Anti-Replay feature was disabled.

Therefore, when encrypting TelePresence over the private WAN and/or MPLS VPN, it is recommended to assign TelePresence to a LLQ and/or to disable Anti-Replay.

TelePresence Branch IPSec VPN Edge QoS Design

As discussed in the previous section, it is not recommended to deploy TelePresence over IPSec VPNs over the Internet, due to the lack of service level guarantees of the Internet in general. But rather, if required due to business reasons, IPSec encryption via the network infrastructure provides an additional security overlay to private WANs or MPLS VPNs for TelePresence calls.

If TelePresence is to be deployed over IPSec VPNs over private WANs or MPLS VPNs, then three additional points should be kept in mind:

- Provision the additional bandwidth required by encryption, according to [Table 6-2](#).
- Provision TelePresence traffic into a LLQ (or a dual-LLQ, along with voice).
- Either maximize or disable Anti-Replay.

The first bullet is straightforward and the second bullet has already been covered (see [Example 6-5](#) and [Example 6-9](#)). The third bullet, however, requires some new commands that we have not yet detailed.

To minimize Anti-Replay errors or to eliminate them completely, Cisco IOS introduced a pair of commands in 12.3(14)T that could either enlarge the Anti-Replay window (to a maximum of 1024 packets) or disable it entirely, the **set security-association replay window-size** and the **set security-association replay disable** commands, respectively.

Let us consider two examples to illustrate these options. In [Example 6-12](#), a dual-LLQ QoS policy (with modified priority bandwidth for the TelePresence class) is applied in conjunction with a native IPSec tunnel (including a maximized Anti-Replay window) to a branch T3 WAN/VPN edge interface.

Example 6-12 Dual-LLQ Design with Native IPsec Tunnel and Maximized Anti-Replay Window

```

!
policy-map WAN-EDGE-IPSEC
  class VOIP
    priority percent 10           ! LLQ for VoIP (example amount of BW)
  class TELEPRESENCE
    priority 15360 256000       ! LLQ for CTS-3000 with IPsec (1080p-Best + aux video)
  class DATA
    ...
!
crypto isakmp policy 1
  encr 3des
  authentication pre-share
  group 2
crypto isakmp key CTS address 192.168.2.10
crypto isakmp keepalive 10
!
!
crypto ipsec transform-set CTS-IPSEC esp-3des esp-md5-hmac
!
crypto map CMAP local-address Serial6/0
crypto map CMAP 10 ipsec-isakmp
  set peer 192.168.2.10
  set security-association replay window-size 1024 ! Maximizes A/R
  set transform-set CTS-IPSEC
  match address BRANCH-TO-CAMPUS
  qos pre-classify
!
!
interface Serial6/0
  description BRANCH-TO-CAMPUS-T3
  ip address 192.168.2.9 255.255.255.252
  tx-ring-limit 10           ! Tunes T3 Tx-Ring
  dsu bandwidth 44210
  framing c-bit
  cablelength 10
  serial restart-delay 0
  crypto map CMAP
  max-reserved-bandwidth 100           ! LLQ/CBWFQ BW Override
  service-policy output WAN-EDGE-IPSEC ! Attaches Dual-LLQ Policy
!
!
ip access-list extended BRANCH-TO-CAMPUS
  permit ip 10.16.0.0 0.0.255.255 10.17.0.0 0.0.255.255
!

```

In [Example 6-13](#), a dual-LLQ QoS policy (with modified priority bandwidth for the TelePresence class) is applied in conjunction with a GRE IPsec tunnel (with a disabled Anti-Replay window) to a branch T3 WAN/VPN edge interface.

Example 6-13 Dual-LLQ Design with GRE IPsec Tunnel and Disabled Anti-Replay Window

```

!
policy-map WAN-EDGE-IPSEC
  class VOIP
    priority percent 10           ! LLQ for VoIP (example amount of BW)
  class TELEPRESENCE
    priority 15360 256000       ! LLQ for CTS-3000 with IPsec
  class DATA
    ...
!

```

```

crypto isakmp policy 1
  encr 3des
  authentication pre-share
  group 2
crypto isakmp key tele address 192.168.2.10
crypto isakmp keepalive 10
!
!
crypto ipsec transform-set CTS-IPSEC esp-3des esp-md5-hmac
!
crypto map CMAP local-address Serial6/0
crypto map CMAP 10 ipsec-isakmp
set peer 192.168.2.10
set security-association replay disable ! Disables Anti-Replay
set transform-set CTS-IPSEC
match address BRANCH-TO-CAMPUS
qos pre-classify
!
!
interface Tunnel0
  ip address 10.18.1.1 255.255.255.252
  tunnel source 192.168.2.9
  tunnel destination 192.168.2.10
!
interface Serial6/0
  description BRANCH-TO-CAMPUS-T3
  ip address 192.168.2.9 255.255.255.252
  tx-ring-limit 10 ! Tunes T3 Tx-Ring
  dsu bandwidth 44210
  framing c-bit
  cablelength 10
  serial restart-delay 0
  crypto map CMAP
  max-reserved-bandwidth 100 ! LLQ/CBWFQ BW Override
  service-policy output WAN-EDGE-IPSEC ! Attaches Dual-LLQ Policy
!
!
ip access-list extended BRANCH-TO-CAMPUS
  permit gre host 192.168.2.9 host 192.168.2.10
!

```

These configurations can be verified with the following command:

- **show policy-map interface**
- **show crypto engine accelerator statistic** *module*

TelePresence Branch MPLS VPN

MPLS VPN architectures are comprised of customer edge (CE) routers, provider-edge (PE) routers, and provider (P) routers. MPLS VPNs provide fully-meshed Layer 3 virtual WAN services to all interconnected CE routers. Let us discuss some of the critical QoS design considerations pertaining to MPLS VPNs and then translate these considerations into configuration examples.



Note

MPLS VPN architectures are defined in RFC 2547 “BGP/MPLS VPNs” at www.ietf.org/rfc/rfc2547.

TelePresence Branch MPLS VPN Edge Considerations

The advent of MPLS VPN service offerings that inherently offer full-mesh connectivity has shifted the QoS administration paradigm. Under traditional hub-and-spoke Layer 2 WAN designs, the enterprise network administrator controlled all the QoS policies by configuring these on the WAN aggregator routers' and branch routers' WAN edges, as previously discussed. However, under a full-mesh topology, it is the service provider's QoS policies on the PE edges routers that ultimately determine how traffic enters a branch and these SP policies may be different from the enterprise's policies on the (unmanaged) CE edges.

Therefore, to ensure end-to-end service levels, enterprise administrators must choose service providers that offer compatible policies to meet their business objectives; furthermore, enterprises must fully understand the SP's QoS policies and map their policies to match in a complementary manner.

First, let us briefly discuss service provider selection based on SLA requirements. As brought out in [Chapter 4, "Quality of Service Design for TelePresence,"](#) the bandwidth and service level requirements of TelePresence (including latency, jitter, and loss requirements) are very high—some are even higher than the SLAs of VoIP. Therefore, to achieve these tight end-to-end SLAs, it is mandatory that the SP be able to guarantee a subset of these SLAs from PE-edge-to-PE-edge.

In the past, to facilitate VoIP deployments over MPLS VPNs, Cisco initiated a Cisco Powered Network (CPN) "IP Multiservice Service Provider" designation that required SPs to offer (independently-verified) PE-to-PE SLA guarantees that would enable enterprise customers to fulfill the end-to-end SLA requirements of VoIP. This initiative was well received by both enterprise customers and SPs, as enterprise customers did not have to do as much research and testing to validate potential SP networks to support their VoIP deployments and SPs, in turn, received a competitive advantage in marketing their networks that could meet VoIP SLAs.

Subsequently, this initiative has similarly been applied to TelePresence. Cisco is in the process of validating various service providers that can meet a subset of the stringent SLAs required by TelePresence, so that enterprise customers can provide end-to-end SLA guarantees for TelePresence.

Second, let us turn our attention to enterprise-to-service provider mapping, which usually involves three main points to consider:

1. The number of enterprise traffic classes versus the number of service provider traffic classes; and if collapsing is required, how to perform this efficiently.
2. Marking or remarking requirements on CE egress to gain admission to the desired SP traffic class; and (optional) remarking requirements on CE ingress to restore enterprise traffic markings for provisioning, accounting, or management purposes.
3. Non-traditional WAN access media, such as sub-line-rate Ethernet access, and the QoS implications these pose.

Let us discuss each of these in turn, beginning with the number of traffic classes.

The number of traffic classes within an enterprise network is a function of its business objectives (as discussed in detail in Chapter 1 of the QoS SRND at www.cisco.com/go/srnd). As an informational guide, RFC 4594 "Configuration Guidelines for DiffServ Service Classes" (www.ietf.org/rfc/rfc4594), outlines up to 12 classes of traffic that may be present within an enterprise. This is not to say that it is mandatory for enterprises to have 12 traffic classes today; but rather that the potential exists in enterprise networks for up to 12 traffic classes and—given the trends in emerging new applications and evolving business objectives—even if enterprises are not deploying 12 class models today, they may need to in the near future. For configuration and testing purposes, we can use this 12-class enterprise model as a worst-case scenario, providing maximum disparity between the number of enterprise traffic classes versus the number of service provider traffic classes. The Cisco-modified 12-class RFC 4594 enterprise model is shown in [Figure 6-9](#).

Figure 6-9 Cisco-Adapted 12-Class RFC 4594-based Enterprise Classification and Marking Model

Application	L3 Classification		IETF
	PHB	DSCP	RFC
Network Control	CS6	48	RFC 2474
VoIP Telephony	EF	46	RFC 3246
Broadcast Video	CS5	40	RFC 2474
Multimedia Conferencing	AF41	34	RFC 2597
Real-Time Interactive/TelePresence	CS4	32	RFC 2474
Multimedia Streaming	AF31	26	RFC 2597
Call Signaling	CS3	24	RFC 2474
Low-Latency/Transactional Data	AF21	18	RFC 2597
Operations/Administration/Management	CS2	16	RFC 2474
High-Throughput/Bulk Data	AF11	10	RFC 2597
Best Effort	DF	0	RFC 2474
Low-Priority/Scavenger Data	CS1	8	RFC 3662

223250

**Note**

Some of the application class names show both the RFC 4594 names as well as the better known, but less wordy, Cisco QoS Baseline application class names. For example, “Low Latency Data,” “High Throughput Data,” and “Low Priority Data” are generally more easily referred to as “Transactional Data,” “Bulk Data,” and “Scavenger,” respectively. Nonetheless, the names can be viewed as synonymous.

Now let us look at service provider class models. At the time of writing, in North America, most service providers offer 3- or 4-class QoS models, although some are planning 6-class models. In EMEA or Asia Pacific, some providers offer even more classes. Rather than presenting models for each number of SP classes, we consider just two, a 4-class and a 6-class model, and the principles applied to these enterprise-to-SP mapping examples can be extended to other traffic class models. These 4-class and 6-class examples are graphically illustrated in [Figure 6-10](#).

Figure 6-10 Example 4-Class and 6-Class MPLS VPN SP QoS Models

4-Class SP Model		6-Class SP Model	
EF CS5	SP-Real-Time (RTP/UDP) 30%	EF CS5	SP-Real-Time (RTP/UDP) 20%
CS6 AF3 CS3	SP-Critical 1 (TCP) 20%	CS4	SP-Critical 1 (TelePresence) 10%
AF2 CS2	SP-Critical 2 (UDP) 20%	CS6 AF3 CS3	SP-Critical 2 (TCP) 20%
DF	SP-Best Effort 30%	AF1 CS1	SP-Scavenger 5%
		DF	SP-Best Effort 25%

223098

Comparing Figure 6-7 to Figure 6-8 highlights the fact that, more often than not, there are generally fewer SP traffic classes than enterprise classes, and thus there are times when more than one enterprise traffic class is assigned to the same SP class. When such collapsing has to be done, it is recommended to **avoid mixing TCP-based applications with UDP-based applications within a single service provider class**. This is due to the behavior of these respective protocols during periods of congestion.

Specifically, due to TCP transmission guarantees and its windowing behavior, TCP transmitters throttle back flows when drops are detected. In contrast, most UDP transmitters are completely oblivious to drops and, therefore, never lower transmission rates because of dropping. When TCP flows are combined with UDP flows within a single service provider class and the class experiences congestion, TCP flows continually lower their transmission rates, potentially giving up their bandwidth to UDP flows that are oblivious to drops. This effect is called TCP starvation/UDP dominance. Even if WRED is enabled on the service provider class, the same behavior would be observed because WRED (for the most part) manages congestion only on TCP-based flows. Granted, it is not always possible to separate TCP-based flows from UDP-based flows, but it is beneficial to be aware of this behavior when making such application-mixing decisions within a single service provider class.

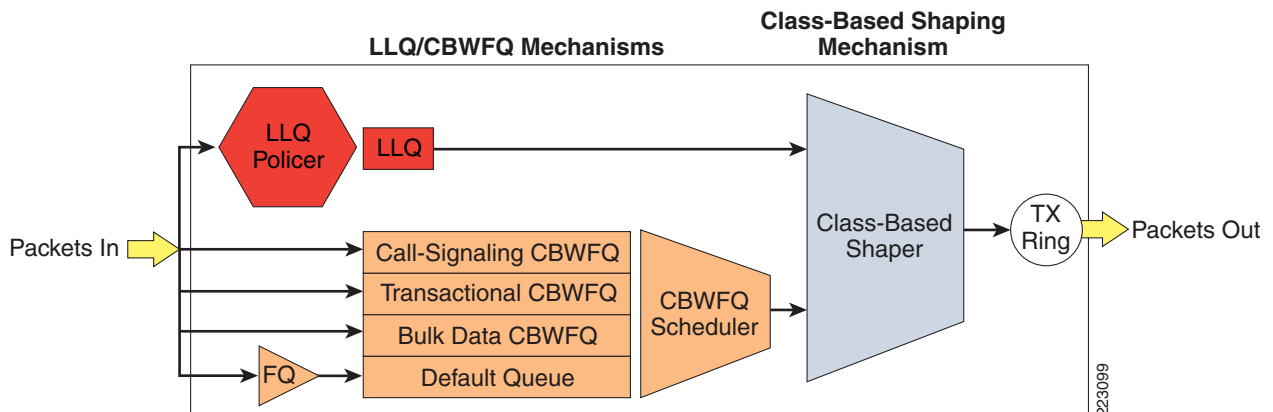
Now let us look at traffic marking and remarking requirements. As can be seen in Figure 6-8, DSCP values serve as the admission criteria per SP class. These DSCP values likely vary from one provider to another, therefore it is important for the enterprise subscriber be fully informed of the DSCP admission criteria for each SP class. At times applications may need to be remarked in order to gain admission to the desired SP class. When such is the case, **marking should be done as the final operations on the (unmanaged) CE egress edge**. Otherwise, if remarking is done at an earlier node, say the campus access edge, then changes to the SP QoS policies or migration to another SP would be much more difficult to manage, as would using multiple SPs for redundancy (each with its own marking scheme).

Also, there may be times when the enterprise has a business requirement to maintain DSCP markings in the branch, perhaps for traffic accounting purposes or for other reasons. In such cases, the enterprise subscriber may choose to make the MPLS VPN appear DSCP-transparent by **restoring enterprise DSCP markings on the CE ingress edge**.

Additionally, each SP class is likely policed on the PE ingress edge. Excess traffic may either be remarked or dropped. Again, it is important for the enterprise subscriber to know exactly how excess traffic is treated on a per-class basis. Understanding SP policing policies is an especially important consideration for the TelePresence class. As we have already discussed in [TelePresence Branch WAN Edge LLQ Policy](#), TelePresence requires 256 KB of committed burst from a policer. **Therefore, it is essential to confirm with the service provider that whatever class TelePresence traffic is assigned to is being policed with at least 256 KB of burst.**

And finally, let us discuss the QoS implications of non-traditional WAN access-media, such as Ethernet. As previously discussed, queuing policies only engage when the physical interface is congested (as is indicated to IOS software by a full Tx-Ring). This means that queuing policies never engage on media that has a contracted sub-line rate of access, whether this media is Frame Relay, ATM, or Ethernet. In such a scenario, **queuing can only be achieved at a sub-line rate by introducing a two-part policy**, sometimes referred to a **Hierarchical QoS (HQoS) policy** or nested QoS policy, **wherein 1) traffic is shaped to the sub-line rate, and 2) traffic is queued according to the LLQ/CBWFQ policies within the sub-line rate.** With such an HQoS policy, it is not the Tx-Ring that signals IOS software to engage LLQ/CBWFQ policies, but rather it is the Class-Based Shaper that triggers software queuing when the shaped rate has been reached. Such an HQoS policy is graphically illustrated in [Figure 6-11](#).

Figure 6-11 Hierarchical QoS Policy—Shaping to a Sub-Line Rate with Queuing within the Shaped Rate

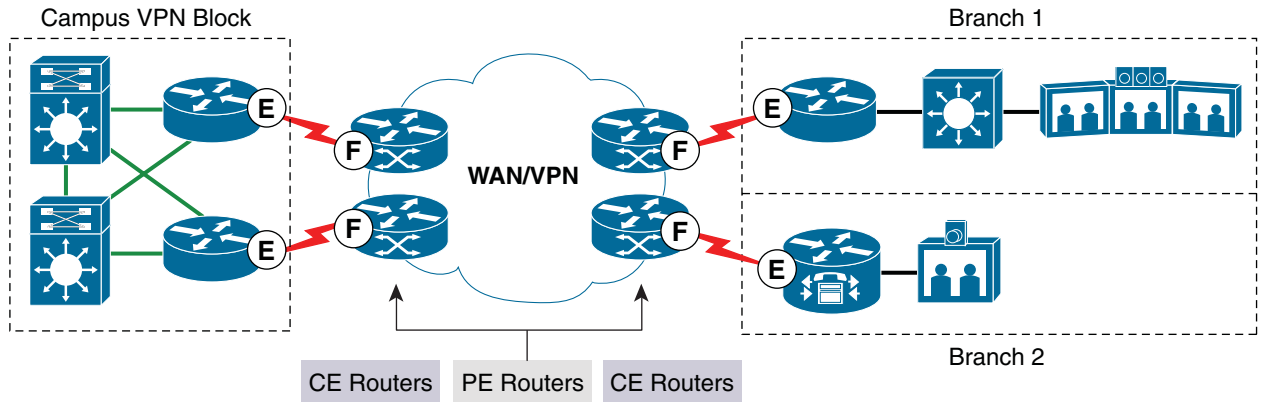


Let us consider a practical example in which an SP offers an enterprise subscriber a GigabitEthernet handoff, but with a (sub-line rate) contract for only 50 Mbps. Normally, queuing policies only engage on this GE interface when the offered traffic rate exceeds 1000 Mbps. However, the enterprise administrator wants to ensure that traffic within the 50 Mbps contracted rate is properly prioritized prior to PE handoff. Therefore, they configure an HQoS policy, such that the interface shapes all traffic to the contracted 50 Mbps rate and attaches a nested queuing policy to the shaping policy, such that traffic is properly prioritized within this 50 Mbps sub-line rate.

The only other consideration an administrator should keep in mind with HQoS policies is their potential performance impact. When performed in IOS software on routers, then these policies generate a marginal CPU load; the actual amount of the load depends on platforms, speeds, policy complexity, traffic rates, and other factors. A rule of thumb, however, is to always keep CPU levels below 75% during normal operating conditions, as this allows some cycles to always be available to process network events. Some platform guidance for HQoS policies are presented, along with detailed configurations, in [TelePresence Branch MPLS VPN QoS Designs](#).

Finally, let us take a look at how all these QoS policies fit together for a TelePresence-enabled branch subscribing to a MPLS VPN, as illustrated in [Figure 6-12](#).

Figure 6-12 Enterprise and Service Provider MPLS VPN QoS Design Recommendations for TelePresence



Enterprise Subscriber (Unmanaged CE Routers):

- (E) Outbound Policies:**
 HQoS Shaper (if required)
 ≤33% of BW { + LLQ for VoIP (EF)
 + LLQ or CBWFQ for TelePresence (CS4)
 + Remark TelePresence (if necessary)
 + CBWFQ for Call-Signaling (CS3)
 + Remark Call-Signaling (if necessary)

- Inbound Policies:**
 Trust for DSCP
 + Restore TelePresence to CS4 (if necessary)
 + Restore TelePresence to CS4 (if necessary)

Service Provider:

- (F) Outbound Policies:**
 + LLQ for Real-Time
 + CBWFQ for Critical Data

- Inbound Policies:**
 Trust for DSCP
 Police on a per-Class Basis

223100

As shown in Figure 6-12, the enterprise subscriber provisions LLQ/CBWFQ policies for VoIP and TelePresence (in conjunction with HQoS sub-line rate shapers, if required) and performs any application-class remarking on the CE egress edges. Optionally, if required, the enterprise may restore their markings on the CE ingress edges for any traffic that required remarking over the MPLS VPN.

In turn, the service provider polices traffic on a per-class basis on their PE ingress edges and provisions LLQ/CBWFQ policies according to their class-models on the PE egress edges. They may also perform QoS and/or MPLS Traffic Engineering within their core; however, such policies are beyond the scope of our enterprise-centric designs.



Note

Due of the explicit ingress policing on PE edges of MPLS VPNs, it cannot be overemphasized that the enterprise subscriber needs a comprehensive Call Admission Control system in place to limit the amount of TelePresence traffic over the MPLS VPN; otherwise the call-quality of **all** TelePresence calls over the MPLS VPN may degrade to the point of unusability.

TelePresence Branch MPLS VPN QoS Designs

Having reviewed the many design considerations for MPLS VPNs, let us now put them into practice by constructing configuration policies to meet the requirements of several specific scenario examples, including a 4-class SP model, a 6-class SP model, and a sub-line rate access example.

TelePresence 4-Class MPLS VPN SP Model QoS Design

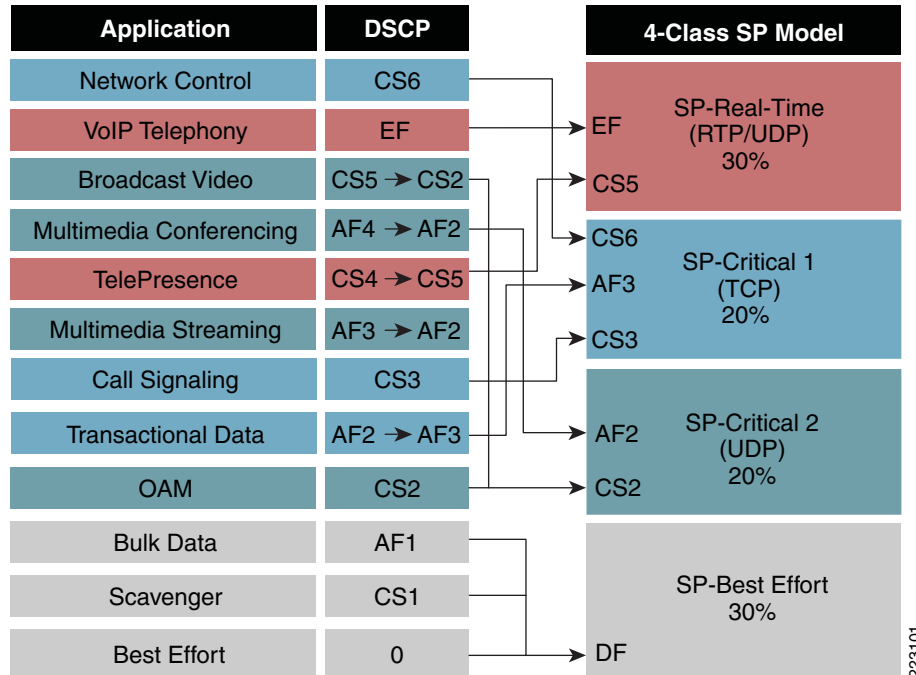
Let us begin by constructing a CE edge policy to support a 4-class SP model example. In this example, since there are so few classes to choose from, TelePresence may need to be combined with another application. **It is highly recommended not to combine TelePresence with any unbounded application** (i.e., an application without any admission control) **within a single SP class**, since this could lead to class congestion, resulting in TelePresence drops (with or without WRED enabled on the SP class), which would ruin TelePresence call quality. Therefore, in such a design two choices exist:

- Assign TelePresence into the SP-Realtime class along with voice.
- Assign TelePresence to a dedicated non-priority SP class.

We consider the option of assigning TelePresence into the SP-Realtime class for this example and then consider the option of assigning it to a non-priority class in the following (6-class SP model) example.

Given the 4-Class SP model illustrated in [Figure 6-10](#), we have a Realtime class, a default Best Effort class, and two additional non-priority traffic classes. In this case, the enterprise administrator may elect to separate TCP-based applications from UDP-based applications by using these two non-priority SP traffic classes. Specifically, if voice and TelePresence are the only applications to be assigned to the SP Realtime class, then Broadcast Video, Multimedia Conferencing, Multimedia Streaming, and Operations/Administration/Management (OAM) traffic (which is largely UDP-based) can all be assigned to the UDP SP-class (SP-Critical 2). This leaves the other non-priority SP class (SP-Critical 1) available for control plane applications, such as Network Control and Call-Signaling, along with TCP-based Transactional Data applications. [Figure 6-13](#) shows the per-class remarking requirements from the CE edge to gain access to the classes within the 4-class SP model, with TelePresence assigned to the SP-Realtime class, along with voice.

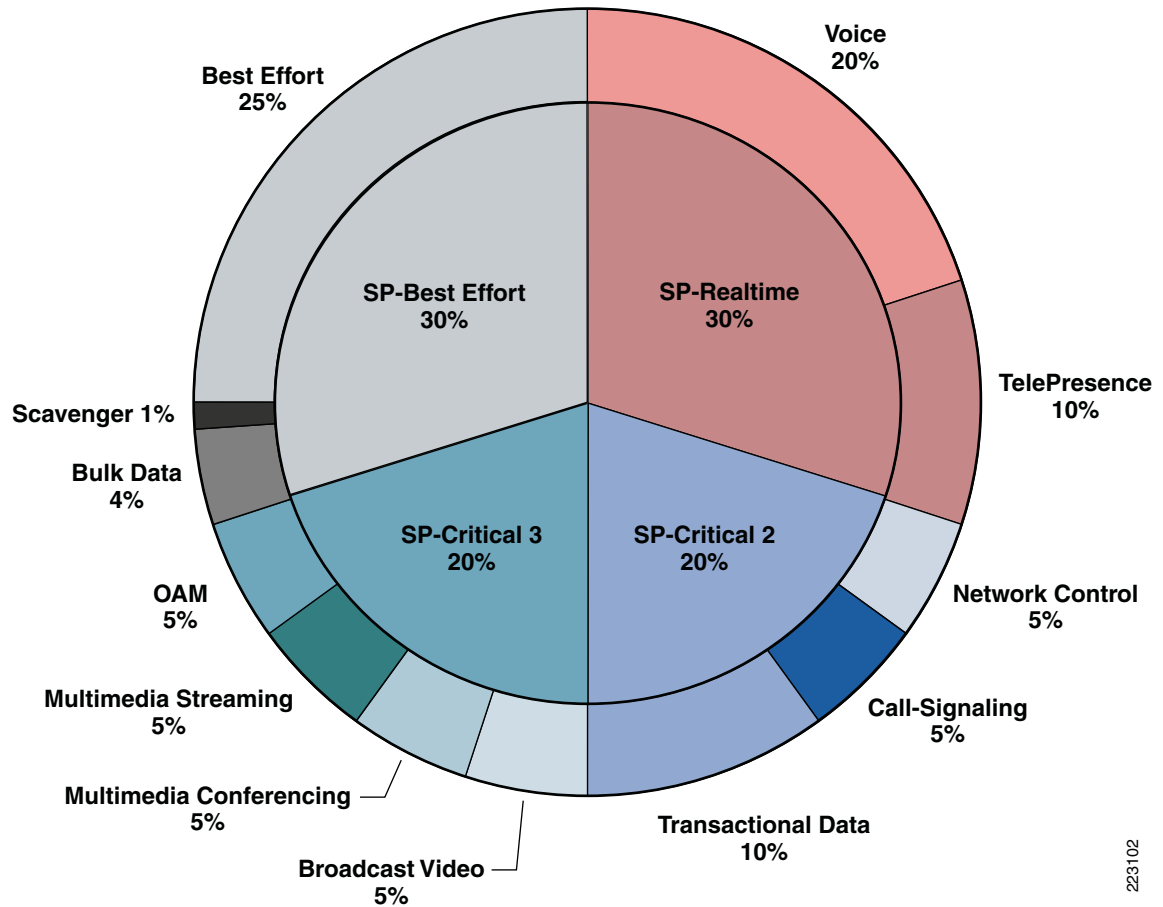
Figure 6-13 Enterprise-to-SP Mapping—4-Class SP Model Example with TelePresence Assigned to the Realtime Class Along with Voice



As shown in [Figure 6-13](#), in this example TelePresence traffic must be remarked on the CE egress edge to CS5 to gain access to the SP's Realtime class. Also, Broadcast Video must be remarked to CS2 to assign it to the UDP SP class (SP-Critical 2). Similarly, Multimedia Conferencing and Multimedia Streaming must be remarked to AF2 to assign these also to the UDP SP class. Correspondingly, Transactional Data traffic must be remarked to AF3 to gain access into the TCP SP class (SP-Critical 1). All other traffic does not require remarking to gain admission to the desired classes; this includes Bulk and Scavenger, as these default to the SP-Best Effort class without any explicit remarking.

Additionally, the relative per-class bandwidth allocations need to be aligned, such that the enterprise CE edge queuing policies are consistent with the SP's PE edge queuing policies to ensure compatible Per-Hop Behaviors (PHBs). Compatible bandwidth allocations are illustrated in [Figure 6-14](#), where the inner pie-chart represents the SP's per-class bandwidth allocations and the outer pie-chart represents the enterprise's per-class bandwidth allocations over an OC3 link.

Figure 6-14 Enterprise-to-SP Bandwidth Allocation C4-Class SP Model Example with TelePresence Assigned to the Realtime Class Along with Voice



The CE egress edge configuration for this policy is shown in [Example 6-14](#).

Example 6-14 Enterprise-to-SP Mapping—4-Class SP Model Example with TelePresence Assigned to the Realtime Class Along with Voice

```

policy-map CE-EDGE-4CLASS-OC3-POS
  class VOICE
    police cir 31000000           ! Voice is policed to 31 Mbps (20%)
    bc 15500                    ! Bc is 15.5 KB
    conform-action transmit      ! Conforming action --> transmit
    exceed-action drop          ! Single-Rate Policing action
    priority                    ! LLQ command for OC3-POS
  class TELEPRESENCE
    police cir 15000000         ! TP is policed to 15 Mbps
    bc 256000                  ! Bc is 256 KB
    conform-action transmit      ! Conforming action --> transmit
    exceed-action drop          ! Single-Rate Policing action
    priority                    ! LLQ command for OC3-POS
    set dscp cs5                ! Remark TelePresence to CS5
  class NETWORK-CONTROL
    bandwidth percent 5        ! CBWFQ for Routing
  class CALL-SIGNALING
    bandwidth percent 5        ! CBWFQ for Call-Signaling
  class TRANSACTIONAL-DATA
    bandwidth percent 10       ! CBWFQ for Transactional Data

```

223102

```

    random-detect dscp-based                ! DSCP-WRED for Transactional Data
    set dscp af31                            ! Remark Transactional Data to AF31
class BROADCAST-VIDEO
    bandwidth percent 5                     ! CBWFQ for Broadcast Video
    set dscp cs2                             ! Remark Broadcast Video to CS2
class MULTIMEDIA-CONFERENCING
    bandwidth percent 5                     ! CBWFQ Video-Conferencing
    random-detect dscp-based                ! DSCP-WRED for Video-Conferencing
    set dscp af21                           ! Remark Video-Conferencing to AF21
class MULTIMEDIA-STREAMING
    bandwidth percent 5                     ! CBWFQ for Streaming-Video
    random-detect dscp-based                ! DSCP-WRED for Streaming-Video
    set dscp af21                           ! Remark Streaming-Video to AF21
class OAM
    bandwidth percent 5                     ! CBWFQ for Network Management
class BULK-DATA
    bandwidth percent 4                     ! CBWFQ for Bulk Data
    random-detect dscp-based                ! DSCP-WRED for Bulk Data
class SCAVENGER
    bandwidth percent 1                     ! Minimum CBWFQ for Scavenger
class class-default
    bandwidth percent 25                    ! CBWFQ for Best Effort
    random-detect                           ! WRED for Best Effort
!
...
...
interface POS3/0/1
    description BRANCH-CE-EDGE-OC3-POS
    ip address 192.168.5.1 255.255.255.252
    clock source internal
    service-policy output CE-EDGE-4CLASS-OC3-POS ! Attaches policy
!

```

Optionally, the original markings for TelePresence, Transactional Data, Broadcast Video, Multimedia Conferencing, and Multimedia Signaling can be restored on the CE ingress edges to make the MPLS VPN appear completely DSCP-transparent to the enterprise, despite the remarking requirements of the service provider. The TelePresence and Transactional Data remarking policies are 1:1 DSCP mappings (one DSCP is changed to another DSCP) and as such are easy to undo with a reversing 1:1 mapping operation. However, the remarking operations performed on Broadcast Video, Multimedia Conferencing, and Multimedia Signaling are 2:1 mappings (two DSCP values are changed to a single DSCP value); specifically Broadcast Video and OAM now share DSCP CS2 and Multimedia Conferencing and Multimedia Signaling share DSCP AF21. These 2:1 DSCP mappings require additional classification policies to identify the discrete applications now sharing a single codepoint. These additional classification policies can include NBAR or access-lists.

In [Example 6-15](#), TelePresence and Transactional Data are restored to their original enterprise-marked DSCP values via a simple 1:1 reverse-mapping. Broadcast Video is separated from OAM by referencing an ACL that identifies the source IP address of the Broadcast Video servers. Multimedia Streaming, in this example, consists of streaming RealAudio and VDO Live stateful protocols, both of which can be identified via NBAR and thus sifted apart from Multimedia Conferencing. An optional DSCP restoration policy for the CE ingress edge is shown in [Example 6-15](#).

Example 6-15 Optional Enterprise DSCP Marking Restoration Policies for CE Ingress Edges

```

class-map match-all SP-TELEPRESENCE
    match dscp cs5                            ! Remark value for TelePresence
class-map match-all SP-TRANSACTIONAL-DATA
    match dscp af31 af32 af33                 ! Remark value(s) for Trans-Data
class-map match-all SP-BROADCAST-VIDEO

```

```

    match dscp cs2                                ! Shared DSCP for Bdcst Video + OAM
    match access-group name BROADCAST-VIDEO-SERVERS! References ACL
class-map match-all SP-MULTIMEDIA-STREAMING-REALAUDIO
    match dscp af21 af22 af32                    ! Shared DSCP for MM-Stream + Conf
    match protocol realaudio                     ! NBAR PDLM for RealAudio
class-map match-all SP-MULTIMEDIA-STREAMING-VDOLIVE
    match dscp af21 af22 af32                    ! Shared DSCP for MM-Stream + Conf
    match protocol vdolive                       ! NBAR PDLM for VDOLive
class-map match-all SP-MULTIMEDIA-CONFERENCING
    match dscp af21 af22 af32                    ! All other AF2 is MM-Conf only
!
policy-map CE-EDGE-IN
  class SP-TELEPRESENCE
    set dscp cs4                                ! Restores original marking for TP
  class SP-TRANSACTIONAL-DATA
    set dscp af21                                ! Restores original marking for TD
  class SP-BROADCAST-VIDEO
    set dscp cs5                                ! Restores original marking for BV
  class SP-MULTIMEDIA-STREAMING-REALAUDIO
    set dscp af31                                ! Restores original marking for MMS
  class SP-MULTIMEDIA-STREAMING-VDOLIVE
    set dscp af31                                ! Restores original marking for MMS
  class SP-MULTIMEDIA-CONFERENCING
    set dscp af41                                ! Restores original marking for MMC
!

interface POS3/0/1
  description BRANCH-CE-EDGE-OC3-POS
  ip address 192.168.5.1 255.255.255.252
  service-policy output CE-EDGE-OC3-POS        ! Attaches egress policy
  service-policy input CE-EDGE-IN              ! Attaches ingress policy
!
...
!
ip access-list extended BROADCAST-VIDEO-SERVERS! Reference ACL
permit ip any 10.200.200.0 0.0.0.255          ! Broadcast Video Server Subnet
!

```

These configurations can be verified with the following command:

- **show policy-map interface**

We can note a few important policy elements in [Example 6-15](#):

- It is important to give the remarked traffic classes unique names from the original enterprise-marked traffic classes, otherwise these interfere or overwrite each other. For example class “TELEPRESENCE” matches on the enterprise marking value for TelePresence (CS4), but class “SP-TELEPRESENCE” matches on the remarked value for TelePresence (CS5).
- Because of the default **match-all** operand on class-maps, we have to have two separate class maps to sift, [traffic marked AF2 and using the realaudio protocol] or [traffic marked Af2 and using the vdolive protocol]. If a single class-map was used, then the **match-all** operand would preclude any match (because the protocol in use cannot be both realaudio and vdolive at the same time; these protocols are unique and represent mutually exclusive criterion). Furthermore, a **match-any** operand on a combined class-map for AF2 **or** realaudio **or** vdolive would not work either, because this would fail the logical requirement that the traffic must be [marked AF2 **and** realaudio] or [marked AF2 **and** vdolive], and thus would result in false-positive matches on Multimedia Conferencing traffic marked AF2.
- It is important to keep in mind that classification logic, like ACL logic, is based on the first-true-match rule. Therefore, the order of classification and sifting must be given careful consideration in order to ensure that traffic marking gets restored properly.

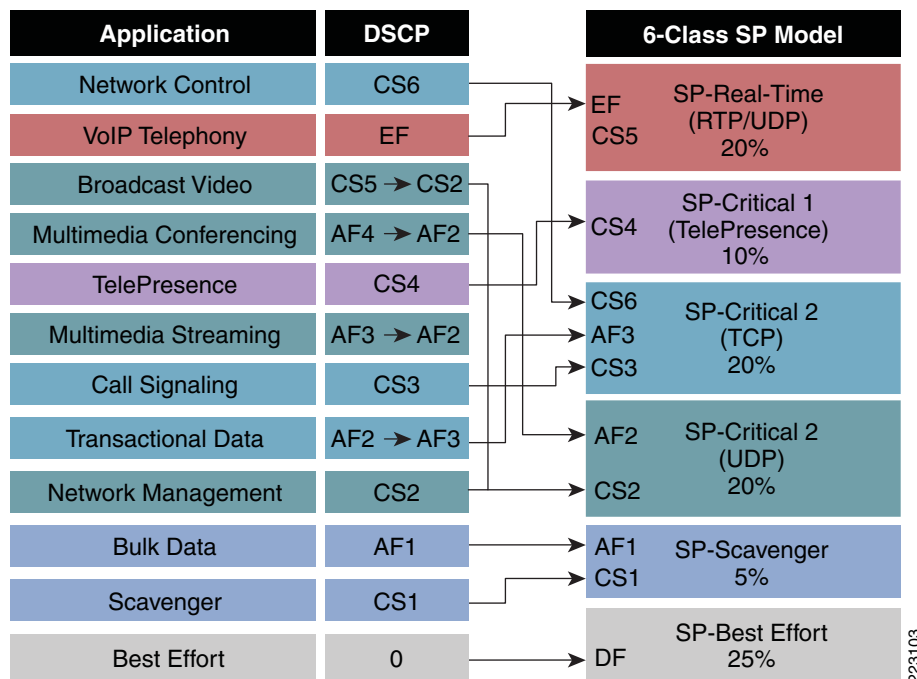
TelePresence 6-Class MPLS VPN SP Model QoS Design

Now let us turn our attention to the 6-Class SP model, also illustrated in [Figure 6-10](#). In this model, we have a Realtime class, a default Best Effort class, a “less-than Best Effort” Scavenger class, and three additional non-priority traffic classes. Furthermore, to illustrate more design options, we assign TelePresence to a non-priority SP-class in this example; but of course TelePresence can also be assigned, in combination with voice, to the SP-Realtime class, as has already been detailed in the previous section.

In this case, the enterprise administrator can dedicate one of the non-priority classes (such as SP-Critical 1) for TelePresence. Again, it bears reiteration that it would not be recommended to assign TelePresence in conjunction with any unbounded application into a single SP class, as the other application could potentially cause the combined class to congest, resulting in TelePresence drops and loss of call-quality.

This leaves two additional non-priority classes, which again allows the administrator to separate TCP-based applications from UDP-based applications. Specifically, Broadcast Video, Multimedia Conferencing, Multimedia Streaming, and Operations/Administration/Management (OAM) traffic can all be assigned to the UDP SP-class (SP-Critical 3). This leaves the other non-priority SP class (SP-Critical 2) available for control plane applications, such as Network Control and Call-Signaling, along with TCP-based Transactional Data applications. [Figure 6-15](#) shows the per-class remarking requirements from the CE edge to gain access to the classes within the 6-class SP model, with TelePresence assigned to a non-priority SP class.

Figure 6-15 Enterprise-to-SP Mapping—6-Class SP Model Example with TelePresence Assigned to a Dedicated, Non-Priority SP-Class

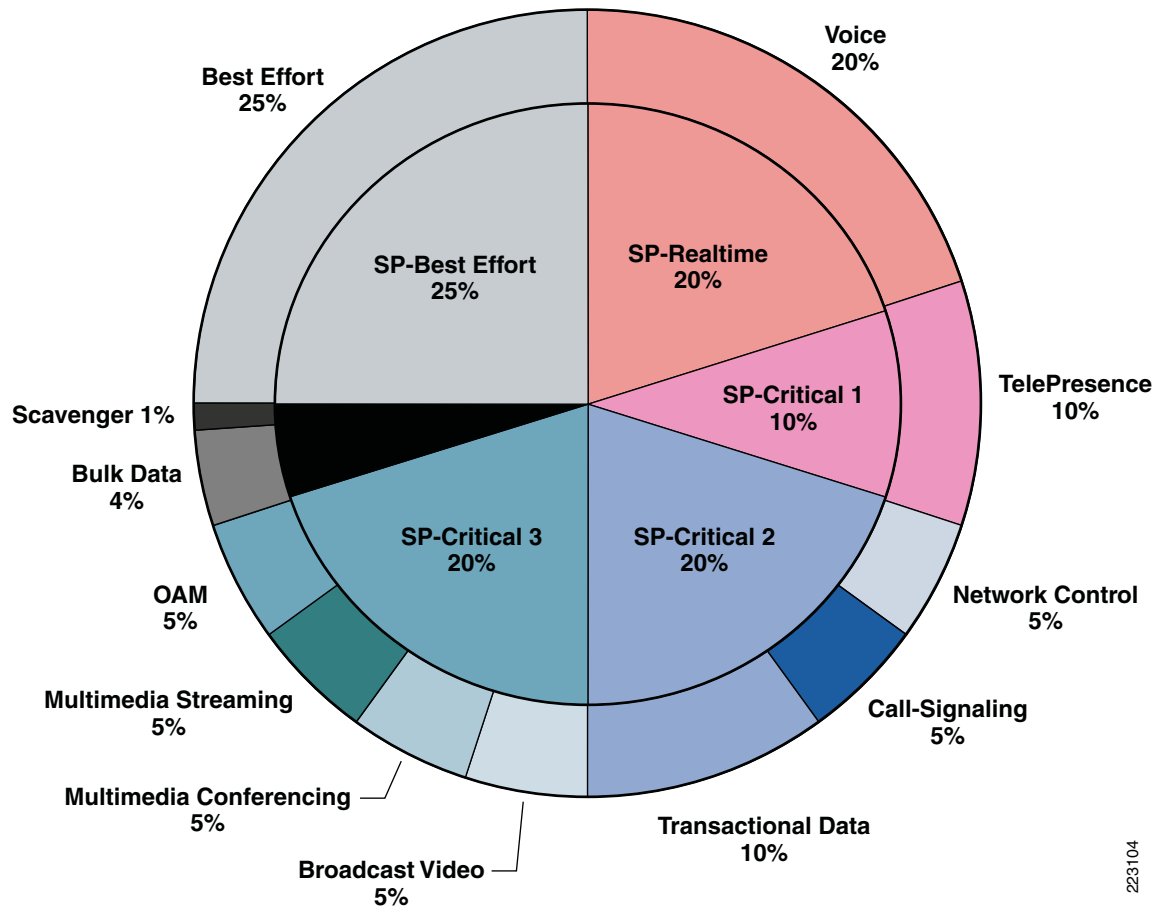


As shown in [Figure 6-15](#), in this second example TelePresence traffic does not need to be remarked to gain access to the dedicated, non-priority SP class to which it is assigned (SP-Critical 1). However as before, Broadcast Video must be remarked to CS2 to assign it to the UDP SP class (SP-Critical 3); Multimedia Conferencing and Multimedia Streaming must be remarked to AF2 to assign these also to the UDP SP class. Correspondingly, Transactional Data traffic must be remarked to AF3 to gain access into the TCP SP class (SP-Critical 2). All other traffic does not require remarking to gain admission to

the desired classes. However, it may be noted that Bulk and Scavenger no longer default to the SP-Best Effort class, but rather now default to the SP-Scavenger class, which is the desired policy to bind these potentially bandwidth-hogging applications.

Additionally, the relative per-class bandwidth allocations again need to be aligned, such that the enterprise CE edge queuing policies are consistent with the SP's PE edge queuing policies. Compatible bandwidth allocations are illustrated in Figure 6-16, where the inner pie-chart represents the SP's per-class bandwidth allocations and the outer pie-chart represents the enterprise's per-class bandwidth allocations over an OC3 link.

Figure 6-16 Enterprise-to-SP Bandwidth Allocation—6-Class SP Model Example with TelePresence Assigned to a Dedicated, Non-Priority SP-Class



The CE egress edge configuration for this policy is shown in Example 6-16.

Example 6-16 Enterprise-to-SP Mapping—6-Class SP Model Example with TelePresence Assigned to a Dedicated, Non-Priority SP-Class

```

policy-map CE-EDGE-6CLASS-OC3-POS
  class VOICE
    police cir 31000000          ! Voice is policed to 31 Mbps (20%)
    bc 15500                   ! Bc is 15.5 KB
    conform-action transmit     ! Conforming action --> transmit
    exceed-action drop         ! Single-Rate Policing action
    priority                   ! LLQ command for OC3-POS
  class TELEPRESENCE

```

```

    bandwidth percent 10                ! CBWFQ for TelePresence
class NETWORK-CONTROL
    bandwidth percent 5                 ! CBWFQ for Routing
class CALL-SIGNALING
    bandwidth percent 5                 ! CBWFQ for Call-Signaling
class TRANSACTIONAL-DATA
    bandwidth percent 10                ! CBWFQ for Transactional Data
    random-detect dscp-based           ! DSCP-WRED for Transactional Data
    set dscp af31                       ! Remark Transactional Data to AF31
class BROADCAST-VIDEO
    bandwidth percent 5                 ! CBWFQ for Broadcast Video
    set dscp cs2                         ! Remark Broadcast Video to CS2
class MULTIMEDIA-CONFERENCING
    bandwidth percent 5                 ! CBWFQ Video-Conferencing
    random-detect dscp-based           ! DSCP-WRED for Video-Conferencing
    set dscp af21                       ! Remark Video-Conferencing to AF21
class MULTIMEDIA-STREAMING
    bandwidth percent 5                 ! CBWFQ for Streaming-Video
    random-detect dscp-based           ! DSCP-WRED for Streaming-Video
    set dscp af21                       ! Remark Streaming-Video to AF21
class OAM
    bandwidth percent 5                 ! CBWFQ for Network Management
class BULK-DATA
    bandwidth percent 4                 ! CBWFQ for Bulk Data
    random-detect dscp-based           ! DSCP-WRED for Bulk Data
class SCAVENGER
    bandwidth percent 1                 ! Minimum CBWFQ for Scavenger
class class-default
    bandwidth percent 25                ! CBWFQ for Best Effort
    random-detect                       ! WRED for Best Effort
!

```

This configuration can be verified with the following command:

- **show policy-map interface**

Optionally, if the original markings for Transactional Data, Broadcast Video, Multimedia Conferencing, and Multimedia Signaling need to be restored, these can be done in a similar manner as demonstrated in [Example 6-15](#), with the exception of not requiring TelePresence traffic to be restored (as it does not get remarked in this 6-class model example).

TelePresence Sub-Line Rate Ethernet Access QoS Designs

As previously discussed, to enforce CE edge queuing policies at sub-line rates, an HQoS policy must be used such that a shaper smooths out traffic to the sub-line rate and forces queuing to occur if this rate is exceeded.

As with policers, Cisco IOS shapers operate on a token-bucket principle, achieving sub-line rates by allowing traffic through in specified bursts (Bc) per sub-second intervals (Tc). As shaping introduces delay to packets above the burst value, it is important to properly size the bursts and intervals to minimize potential shaping jitter. For example, as previously discussed, 1080p sends 30 frames of video per second or, phrased differently, a frame's worth of information every 33 ms. However, if the shaping interval is set too low, say to 5 ms, then a frame's worth of information may be delayed over 3-5 shaping intervals (depending on the amount of frame information). Extensive lab testing has shown that configuring a shaping interval of 20 ms has resulted in the most consistent and minimal jitter values to support a CTS-3000 call.

The interval parameter cannot be set directly, but is set indirectly by explicitly configuring the burst parameter. The relationship between the interval, burst, and shaped rate is given as:

$$Tc = Bc / \text{Shaped Rate}$$

Or:

$$Bc = \text{Shaped Rate} * Tc$$

For example, on a FE or GE interface configured to support a sub-line rate of 50 Mbps, a burst value of 1 megabit (50 Mbps * 20 ms) would result in an optimal shaping interval for TelePresence.



Note

For Cisco IOS Shapers, like the Class-Based Shaper, the burst is expressed in bits (not in Bytes, as is the case with policers).

Translating this into an HQoS policy yields the following configuration, as shown in [Example 6-17](#).

Example 6-17 HQoS Policy to Queue and Shape TelePresence Traffic to a 50 Mbps Sub-Line Rate Over a GigabitEthernet Interface

```

policy-map CE-EDGE                                ! CE Edge queuing policy
  class VOICE
  ...
  class TELEPRESENCE
  ...                                             ! Either a dual-LLQ or CBWFQ policy for TP
  ...
  !
policy-map HQoS-50MBPS                            ! CE Edge HQoS Shaping policy
  class class-default
    shape average 50000000 1000000              ! Rate=50Mbps; Bc=1Mb, Tc=20ms
    service-policy CE-EDGE                      ! Forces queuing at sub-line rate
  !
  ...
  !
interface GigabitEthernet0/1
  description CE-EDGE-GE
  ip address 192.168.1.50 255.255.255.252
  no ip redirects
  no ip proxy-arp
  duplex auto
  speed auto
  media-type rj45
  negotiation auto
  service-policy output HQoS-50MBPS            ! Attaches HQoS policy to GE int
  !

```

This configuration can be verified with the following command:

- **show policy-map interface**

The final point of consideration for this chapter is the performance impact of HQoS policies on various platforms. As previously discussed, when QoS is performed in hardware, such as on Catalyst switches, then there is no performance impact of QoS policies on the CPU. However, when QoS policies are performed in software, then there is a performance impact that depends on several factors, such as the platform, speed, traffic mix, QoS policy, etc.

At speeds up to 150 Mbps, Cisco IOS routers, like the 3800 series ISR or the Cisco 7200-VXR, may be used to enforce HQoS policies on Ethernet-based access-media. Before we look at the performance impact on these platforms, it bears mentioning that it is not primarily the speed that affects the CPU performance, but rather the Packets-per-second (PPS) rate.

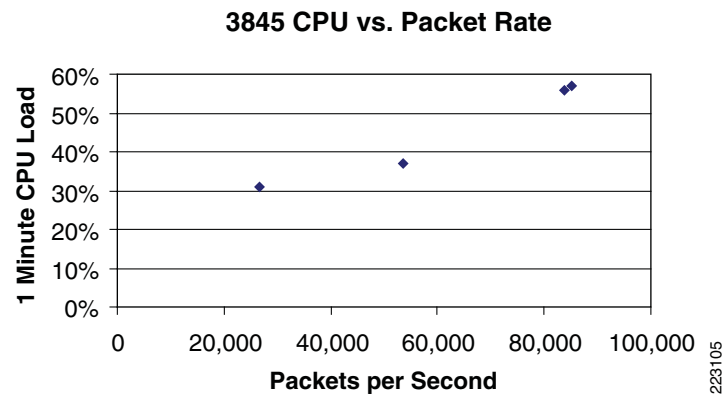
That being said, [Table 6-3](#) shows the performance of a Cisco 3845 router enforcing HQoS policies at rates ranging from 50 Mbps through 150 Mbps.

Table 6-3 Cisco 3845 Platform Performance of HQoS Policies at 50 Mbps Through 150 Mbps Traffic Rates (26KPPS Through 84 KPPS)

Bidirectional Target Data Load	Actual Data Load	CPU	PPS	Video Quality
50 Mbps	50 Mbps Out / 48 Mbps In	31%	26,568	Near Perfect
100 Mbps	87 Mbps Out / 88 Mbps In	37%	53,633	Near Perfect
150 Mbps	145 Mbps Out / 147 Mbps In	57%	85,214	Near Perfect
150 Mbps	145 Mbps Out / 140 Mbps In	56%	83,879	Near Perfect

The corresponding performance graph for [Table 6-3](#) is shown in [Figure 6-17](#).

Figure 6-17 Cisco 3845 Platform Performance of HQoS Policies at 50 Mbps Through 150 Mbps Traffic Rates (26KPPS Through 84 KPPS)



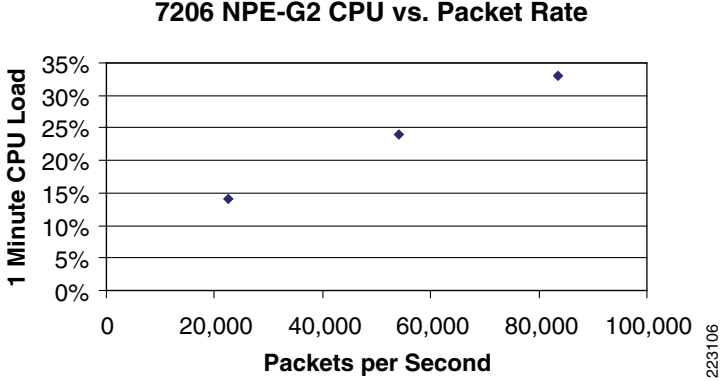
Additionally, [Table 6-4](#) shows the performance of a Cisco 7200VXR router with a Network Processing Engine (NPE) G2 enforcing HQoS policies at rates ranging from 50 Mbps through 150 Mbps.

Table 6-4 Cisco 7200VXR with NPE-G2 Platform Performance of HQoS Policies at 50 Mbps Through 150 Mbps Traffic Rates (22KPPS Through 84 KPPS)

Bidirectional Target Data Load	Actual Data Load	CPU	PPS	Video Quality
50 Mbps	50 Mbps Out / 45 Mbps In	14%	22,634	Near Perfect
100 Mbps	96 Mbps Out / 96 Mbps In	24%	54,011	Near Perfect
150 Mbps	142 Mbps Out / 147 Mbps In	33%	83,631	Near Perfect

The corresponding performance graph for [Table 6-4](#) is shown in [Figure 6-18](#).

Figure 6-18 Cisco 7200VXR with NPE-G2 Platform Performance of HQoS Policies at 50 Mbps Through 150 Mbps Traffic Rates (22KPPS Through 84 KPPS)



Both of these platforms are able to support HQoS policies for TelePresence at these speeds (50 Mbps through 150 Mbps). The goal, however, is to keep CPU levels below 75% during normal conditions, so that the router always has some cycles available to process network events.

For higher speeds, HQoS policies should be performed in hardware (such as on the Catalyst 3750-Metro switch with Enhanced Services modules) or on a hybrid hardware/software platform like the Cisco 7600 SIP/SPA combination.



CHAPTER 7

Call Processing Overview

Overview

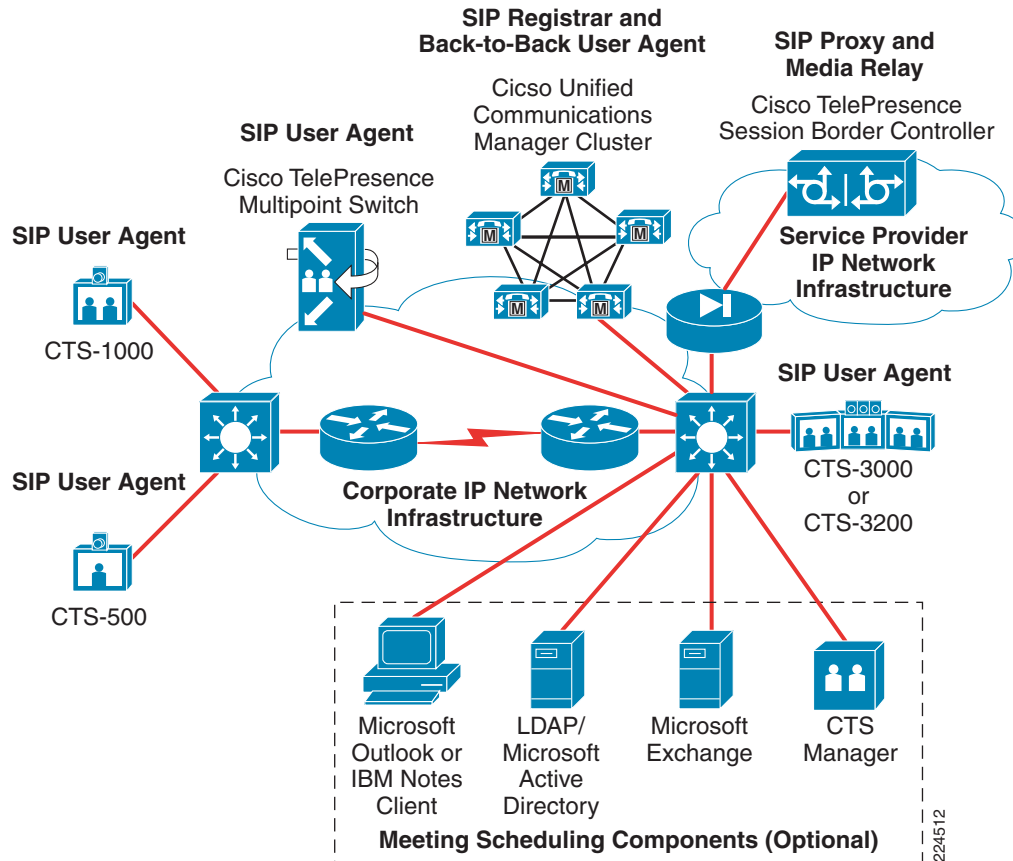
This chapter discusses the Session Initiation Protocol (SIP) and call processing design for Cisco TelePresence, including:

- How the Cisco TelePresence suite of virtual meeting solutions integrates with Cisco Unified Communications Manager (CUCM)
- CUCM software version requirements
- Current CUCM cluster design recommendations
- How the Cisco TelePresence Codecs use Session Initiation Protocol (SIP) and how they register using a shared line appearance with the Cisco Unified 7975G IP phone
- How Cisco TelePresence multipoint resources, such as the Cisco TelePresence Multipoint Switch (CTMS), are configured as a SIP trunk to CUCM and how multipoint calls are routed

Call Processing Components

[Figure 7-1](#) shows the components involved in point-to-point and multipoint TelePresence meetings.

Figure 7-1 Cisco TelePresence Solution Components



These components consist of:

- Two or more Cisco TelePresence systems (any combination of CTS-300s, CTS-3200s, CTS-1000s, or CTS-500s), each with a Cisco Unified 7975G IP phone (not shown in Figure 7-1) which functions as the user interface for launching, controlling, and concluding the meeting

- One CUCM Cluster

TelePresence requires CUCM version 5.1.1 or higher, with version 5.1.2 recommended for support of the Auto Collaborate feature.

- One or more Cisco TelePresence Multipoint Switches (required for multipoint TelePresence meetings)
- IP network infrastructure over which the signaling, video, and audio media are transported
- A Cisco TelePresence Session Border Controller (SBC) typically used for inter-Enterprise TelePresence calls
- Meeting scheduling components (optional):
 - Microsoft Exchange 2003 server
 - Microsoft Active Directory 2000 or 2003 server
 - Microsoft Outlook client
 - Cisco TelePresence Manager (CTS-MAN)

These components are only required for scheduled TelePresence meetings. Ad hoc and permanent TelePresence meetings do not require them.

TelePresence Endpoint Interface to CUCM (Line-Side SIP)

CUCM is the core call processing software for the Cisco TelePresence solution as well as all other Cisco IP telephony devices. CUCM functions as both a SIP registrar and Back to Back User Agent (B2BUA). TelePresence Codex and 7975G IP phones use SIP for call signaling and control, functioning as SIP user agents which register with a CUCM cluster. Cisco TelePresence Systems use TCP for their SIP signaling to/from CUCM. It should be noted that TelePresence devices are currently not supported by the Survivable Remote Site Telephony (SRST) feature of Cisco router platforms, which is often used to provide resiliency in CUCM deployments with remote sites.

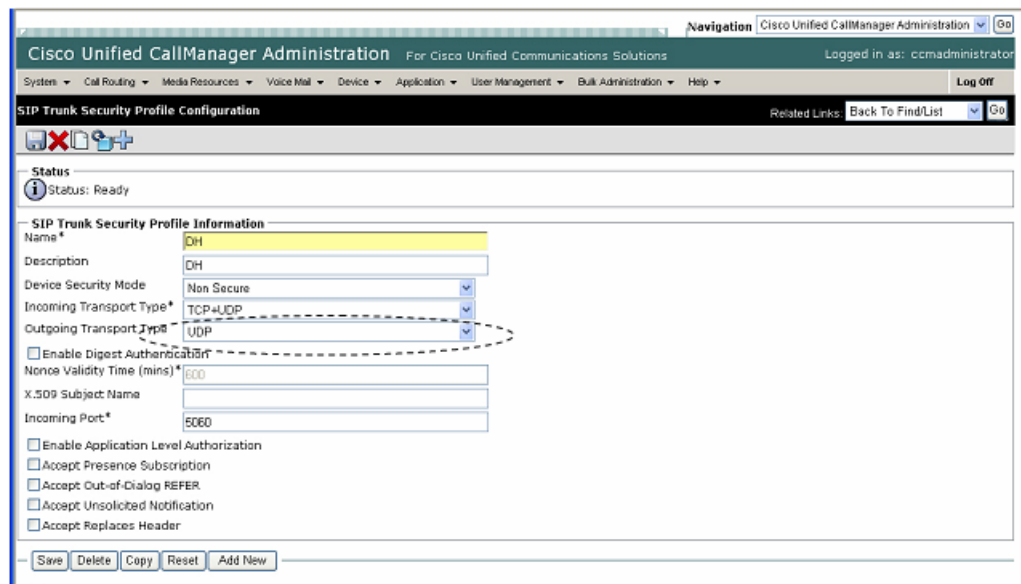
The following sections provide an overview of how TelePresence components register with CUCM, initiate a meeting, and then conclude a meeting.

TelePresence Multipoint Switch Interface to CUCM (Trunk-Side SIP)

The Cisco TelePresence Multipoint Switch (CTMS) multipoint solution connects to CUCM by way of a SIP Trunk. SIP trunks do not use the SIP REGISTER method, and thus for trunks CUCM functions solely as a Back-to-Back User Agent (B2BUA). Route Pattern(s) are configured to route multipoint calls to the SIP trunk(s) of the multipoint switch(es). Prior to software version 1.1, CTMS supported only UDP for SIP signaling to/from CUCM. As of software version 1.1, CTMS also supports TCP.

Therefore the outgoing transport type on the CUCM SIP Trunk Security Profile Configuration must be set for UDP for CTMS. This is shown in [Figure 7-2](#).

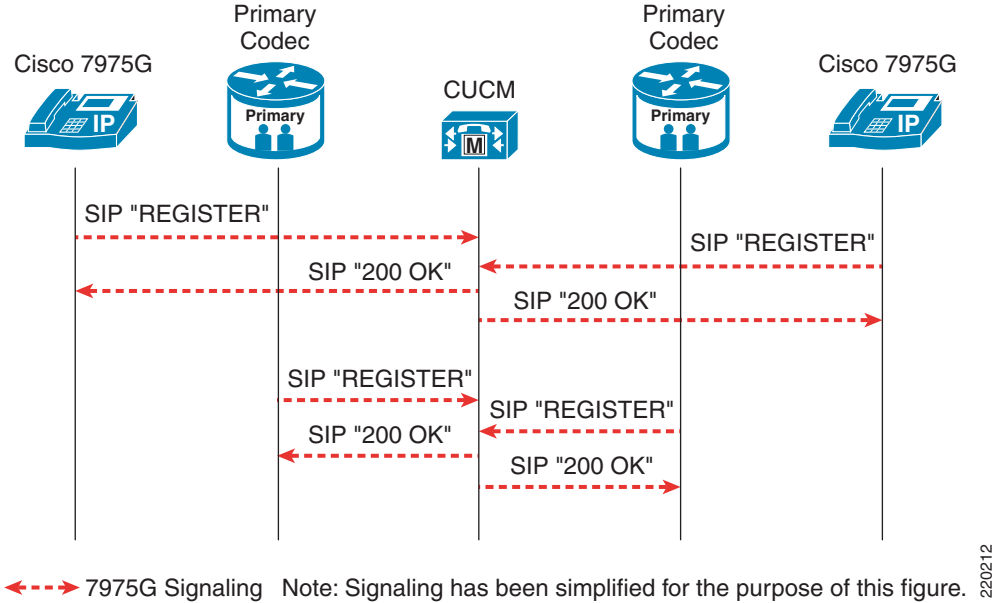
Figure 7-2 CUCM SIP Trunk Security Profile Configuration



TelePresence Endpoint Device Registration

For the initial release of the TelePresence solution, it is recommended that all TelePresence devices in a deployment register to a single CUCM cluster. Although TelePresence devices can be registered across multiple CUCM clusters, Cisco TelePresence Manager (CTS-MAN), which performs meeting scheduling, can only support a single CUCM cluster in the current release. The 7975G IP phones which function as the user interface for the TelePresence solution also register with CUCM, sharing the same dial extension as the TelePresence Codecs. Figure 7-3 shows an example of the high-level data flows in the registration process.

Figure 7-3 Cisco TelePresence Device Registration

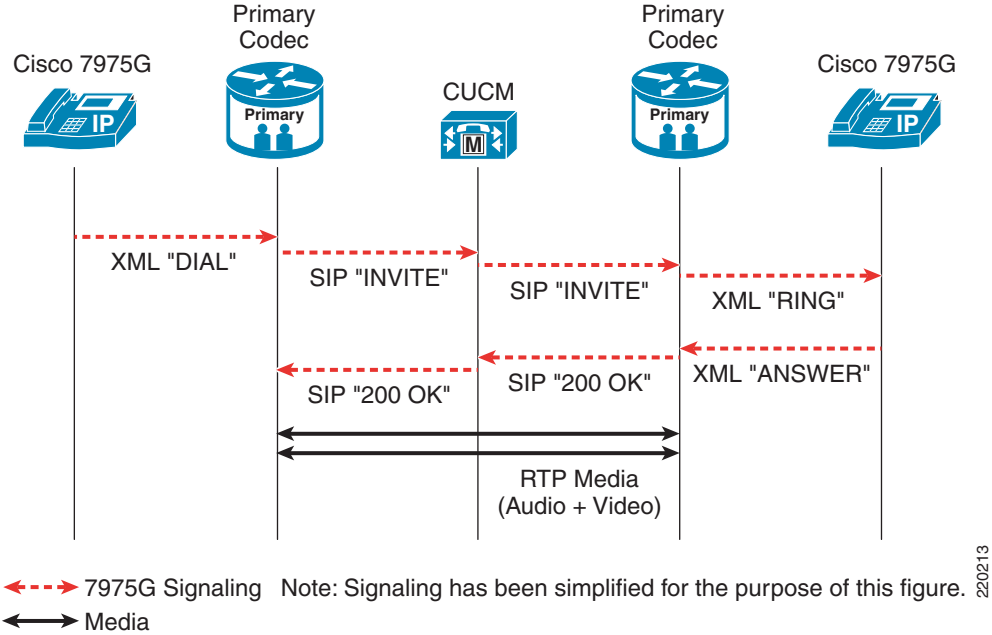


By default CUCM listens on TCP and UDP port 5060 for SIP-related signaling. Cisco TelePresence Systems and Cisco 7975G IP Phones use TCP and hence connect to CUCM on TCP port 5060. The contact header within the SIP REGISTER provides the IP address, transport protocol, port number, and the dial extension for CUCM to reach the TelePresence Codecs and 7975 IP phones.

Call Setup

Once registration is complete, meetings may be established between any two Cisco TelePresence systems or between any TelePresence System and a multipoint switch. Figure 7-4 shows a high-level overview of the call establishment signaling between TelePresence Codecs, their associated 7975G IP phones, and the CUCM cluster.

Figure 7-4 Point-to-Point Cisco TelePresence Call Setup



To make the SIP signaling easier to understand, it has been greatly simplified in Figure 7-4. SIP SUBSCRIBE and NOTIFY messages have been removed from the call flow. These messages are used primarily to update the 7975G IP phones and TelePresence Codecs regarding the status of the call. Finally, HTTP messages between TelePresence Codecs and the Cisco TelePresence Manager have also been removed. These messages inform the Cisco TelePresence Manager of the beginning and ending of a TelePresence meeting.

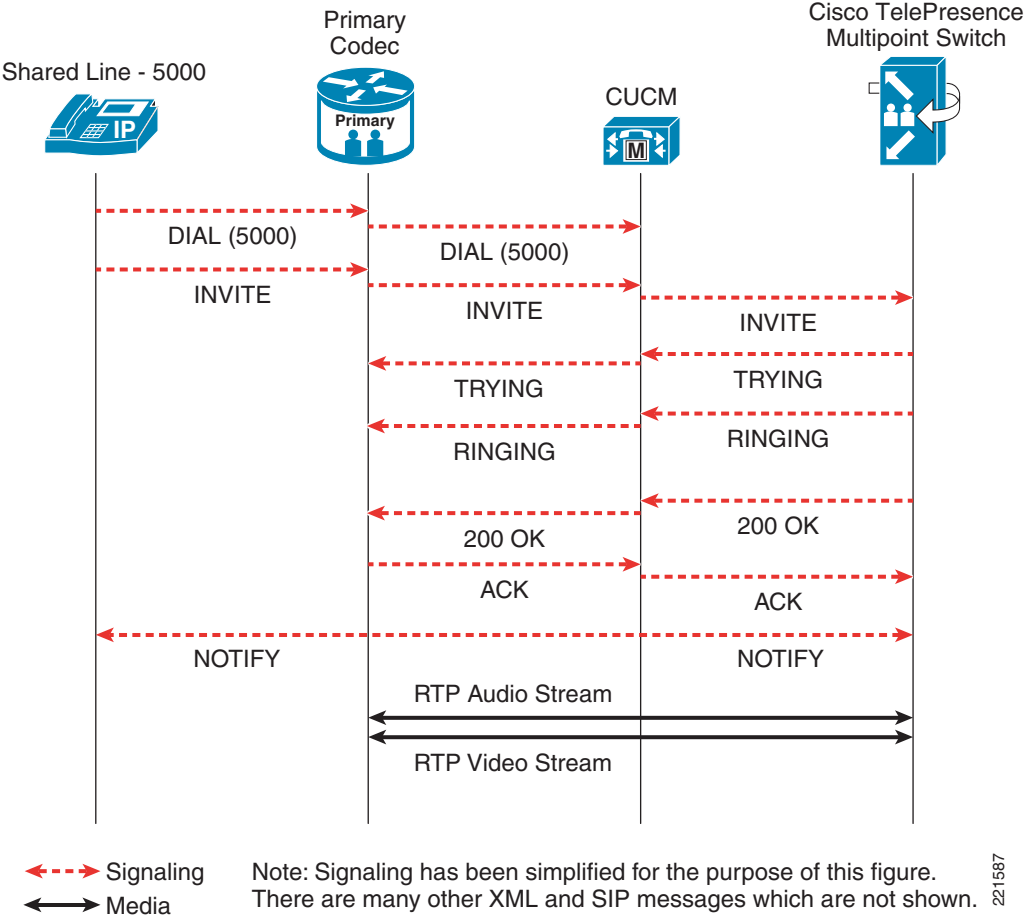
Call setup is initiated when the end user enters or selects, via the touch-screen user interface of the 7975G IP phone, the remote TelePresence location to which he or she wishes to establish a meeting. This causes the 7975G IP phone to generate an XML message to the TelePresence Codec. The XML message instructs the TelePresence Codec to generate a SIP INVITE, which is sent to the CUCM cluster. Within the initial SIP INVITE, the TelePresence Codec uses the Session Description Protocol (SDP). SDP, discussed in IETF RFC 2327, allows two endpoints which are configured for different audio or video modes to negotiate a common set of media parameters for the call. This is accomplished primarily through the use of the media (m=...), attribute (a=...), and bandwidth (b=...) lines. The quality parameter within the TelePresence device configuration in CUCM determines what media capabilities are offered in the initial SDP.

Upon receiving an INVITE from one TelePresence System and determining the destination endpoint (based on the number dialed), CUCM generates a new SIP INVITE to the remote TelePresence Codec. Upon receipt of the SIP INVITE, the TelePresence Codec informs the 7975G IP phone of the incoming call via an XML message. The end user at the remote location accepts the incoming call via the touch-screen user interface of the 7975G phone. This causes a final XML message to be sent to the remote TelePresence Codec, informing it to answer the call. After that, the audio and video media streams begin. Optionally, the TelePresence codec may be configured (in CUCM) to automatically answer all incoming calls, in which case the XML message sequence to/from the phone is skipped and the call is answered immediately. Incidentally, it should be noted that since the same dial extension is shared between the remote TelePresence Codec and the remote 7975G IP phone which functions as its user interface, CUCM generates the new SIP INVITE message to both remote devices. This allows the user to answer the call using the **handset** of the IP Phone (in which case the call is established as an audio-only call). Under normal conditions though, the TelePresence Codec is the one to answer the call and the SIP INVITE to the 7975G IP Phone is canceled.

CUCM acts as a back-to-back user agent (B2BUA), processing requests as a user agent server (UAS) and generating requests as a user agent client (UAC). Unlike a proxy server, CUCM maintains dialog state and participates in all requests sent on the dialogs it establishes. Since CUCM functions as a B2BUA, it sees the SDP information regarding the media capabilities of both sides of the TelePresence call. It determines what audio and video parameters are used for the meeting based on the parameters that are common to both TelePresence devices and what is allowed via the configuration within CUCM. The configuration parameters for the allowed audio and video rates are based on two things: the Quality Setting for each TelePresence System (e.g. 1080p-Best, 1080p-Better, 1080p-Good, 720p-Best, 720p-Better and 720p-Good) and the region settings of the device pool to which the TelePresence devices belong. This allows CUCM to set up a call between two TelePresence devices which are configured for different video modes. For example, if one TelePresence device is configured for 1080p-Best while another is configured for 720p-Good, CUCM specifies 720p in the outgoing SIP message to the 1080p system, thereby negotiating the call down to 720p in both directions.

Multipoint calls are no different than point-to-point calls in that each TelePresence System dials the number of the multipoint switch in a point-to-point fashion. In other words, a multipoint call is nothing more than several point-to-point calls all landing on the same destination device (the multipoint switch). The differences are that instead of matching the dialed number to a Directory Number assigned to a registered endpoint, CUCM matches the dialed number to a Route Pattern assigned to a SIP trunk. The signaling and media negotiation sequences are otherwise the same.

Figure 7-5 Multipoint Cisco TelePresence Call Setup

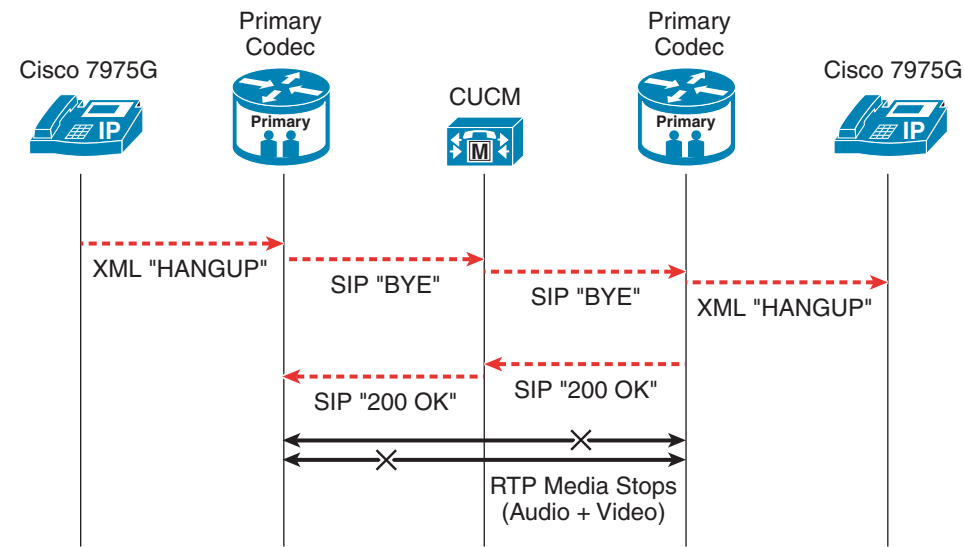


TelePresence utilizes a single AAC-LD over RTP audio stream and a single H.264 over RTP video stream in each direction, for a total of four RTP media streams per bi-directional point-to-point TelePresence meeting. This holds regardless of the model of Cisco TelePresence system device. With CTS-3000 devices, the video streams from the multiple cameras are multiplexed into a single RTP stream. Likewise, the audio streams are multiplexed into a single audio stream. The auxiliary video and audio streams are also multiplexed into these streams.

Call Teardown

Figure 7-6 shows a high-level overview of the call termination signaling between TelePresence Codecs, the 7975G IP phones which function as their user interfaces, and the CUCM cluster.

Figure 7-6 Cisco TelePresence Call Termination



←- - - - -> 7975G Signaling Note: Signaling has been simplified for the purpose of this figure. There are many other XML and SIP messages which are not shown.

←- - - - -> Media

220214

To make the SIP signaling easier to understand, it has again been greatly simplified in Figure 7-6. Call termination begins when the end user at one end of a TelePresence meeting uses the touch-screen user interface of the 7975G IP phone to end the meeting. This causes the 7975G IP phone to send an XML message to the TelePresence Codec, instructing it to hang up the call by generating a SIP BYE message. The SIP BYE message is sent to CUCM, which then generates a new SIP BYE message to the remote TelePresence Codec. The remote TelePresence Codec informs the 7975G phone at the remote site that the call is terminating. Upon receipt of the SIP 200 OK messages from the TelePresence Codecs, the audio and video media streams stop.

Since CUCM functions as a B2BUA which maintains state of all SIP calls initiated and terminated through it, it can capture call detail records of when TelePresence meetings start and stop. This may be necessary for management systems and for billing charges for TelePresence meetings back to individual departments.

Firewall and NAT Considerations

TelePresence embeds the audio and video media endpoint addresses within the SIP call signaling messages. This has implications for firewalls and network address translation. For a firewall to determine the IP addresses and ports to dynamically open to allow the audio and video media through, the firewall may need to monitor the SIP signaling flow. Also, any IP address translation within the network may require special handling, since the addressing received by the remote TelePresence device may not represent a routable IP address to the routers and Layer 3 switches at the remote site. Firewalling within an intra-enterprise TelePresence deployment is discussed in [Chapter 13, “Internal Firewall Deployments with Cisco TelePresence.”](#)



CHAPTER 8

Capacity Planning and Call Admission Control

Overview

The Cisco TelePresence suite of virtual meeting solutions supports three different types of meetings which may be implemented within the Intra-Enterprise Deployment Model:

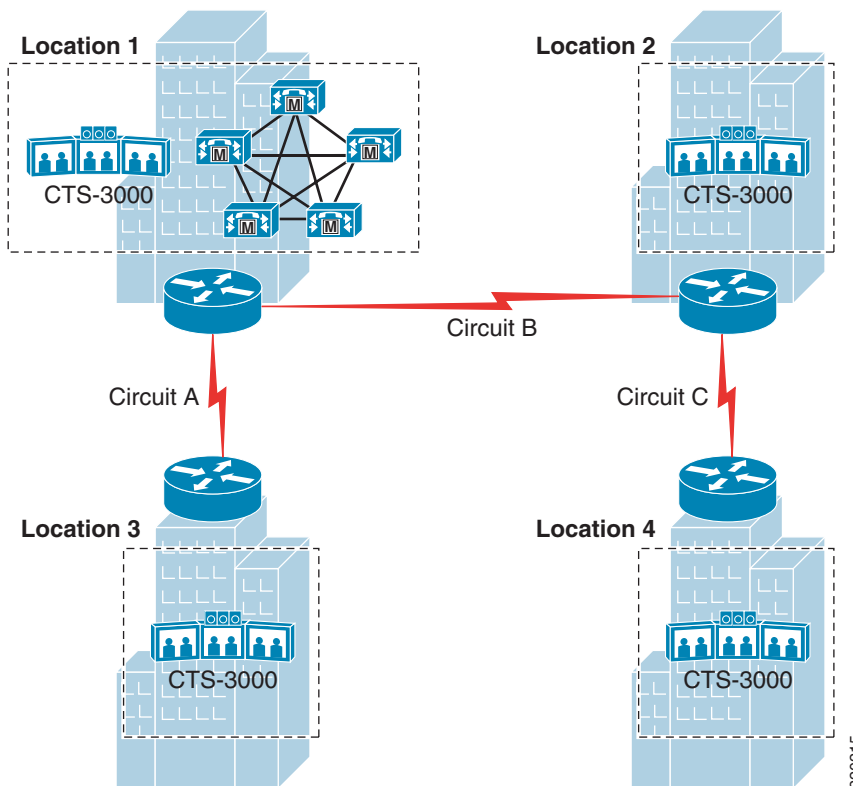
- **Ad hoc meetings**—An end-user simply dials the extension of the Cisco TelePresence system at the other end through the 7975G IP phone that functions as the user interface to the Cisco TelePresence system. There is no scheduling involved. Note that in a multipoint TelePresence meeting, static meetings function in this manner. Ad hoc multipoint meetings must be set up by a CTMS meeting scheduler or administrator.
- **Permanent meetings**—Remain up at all times. An example of a permanent TelePresence meeting is the use of a remote receptionist. Also, in scenarios where there are only two TelePresence systems deployed and they are heavily used, it may be desirable to simply leave the meeting up continuously.
- **Scheduled meetings**—Scheduled in advance of the meeting through the company's groupware application (e.g., Microsoft Exchange/Outlook).

With the current release of the Cisco TelePresence Solution, there is no automated mechanism for reserving network bandwidth or performing call-by-call Call Admission Control (CAC). Therefore, if the number of TelePresence rooms deployed at a given site exceeds the bandwidth available to/from that site, it is possible that too many TelePresence meetings could occur simultaneously and QoS policies in the network will begin dropping TelePresence packets, resulting in poor audio and video quality for all calls traversing that network link. Existing CAC techniques, which are Locations-based CAC or Resource ReserVation Protocol (RSVP), both of which are administered by Cisco Unified Communications Manager (CUCM), are not recommended or supported for Cisco TelePresence. Therefore, the current recommendation is to use manual capacity planning to provide sufficient bandwidth to support all possible TelePresence meetings simultaneously occurring across the network infrastructure. However, due to the limitations of this approach, more advanced CAC mechanisms for TelePresence are being developed and evaluated.

Manual Capacity Planning

Manual capacity planning relies on having sufficient bandwidth within the network to support all possible TelePresence meetings occurring simultaneously and so guarantee 100% call completion. Since all TelePresence meetings are always allowed onto the network, this technique may also be referred to as having no CAC. The physical topology of the network infrastructure impacts how much and where bandwidth needs to be provisioned. [Figure 8-1](#) shows an example of this technique with four locations in a partially-meshed network topology.

Figure 8-1 Bandwidth Provisioning Example



One technique for determining the amount of bandwidth required across each circuit is to simply list all possible combinations of simultaneous TelePresence meetings between locations and the number of meetings each circuit must handle, as shown in [Table 8-1](#).

Table 8-1 Circuit Requirements Example

Meetings Between Locations	Circuit Requirements
Location 1 to Location 2 and Location 3 to Location 4	Circuit A-1 Meeting Circuit B-2 Meetings Circuit C-1 Meeting
Location 1 to Location 3 and Location 2 to Location 4	Circuit A-1 Meeting Circuit B-0 Meetings Circuit C-1 Meeting
Location 1 to Location 4 and Location 2 to Location 3	Circuit A-1 Meeting Circuit B-2 Meetings Circuit C-1 Meetings

However, for the simple network topology shown in [Figure 8-1](#), it is obvious by simply visualizing the network that circuit B must be provisioned with sufficient bandwidth to support two TelePresence meetings, while circuits A and C must be provisioned with sufficient bandwidth to support one TelePresence meeting. Note that for converged networks, this bandwidth is in addition to any other VoIP or video applications, as well as all data traffic. Also, for simplicity, all the devices in [Figure 8-1](#) are

220215

shown as CTS-3000 units. The amount of bandwidth required per Cisco TelePresence meeting depends on the Cisco TelePresence system models (CTS-500, CTS-1000, CTS-3000, or CTS-3200) involved in the call and the video mode (1080p or 720p) which the units are configured to use. The network administrator must take these issues into consideration when determining the amount of bandwidth that must be provisioned to support TelePresence meetings across the network infrastructure. See [Table 4-1](#) in [Chapter 4, “Quality of Service Design for TelePresence”](#) for a detailed list of bandwidth requirements per system type.

The design objective of 100% call completion for all scheduled, ad hoc, and permanent TelePresence meetings is feasible and desirable for current deployments consisting of dozens to hundreds to systems. However, as the number of TelePresence endpoints deployed increases into the hundreds or even thousands, the amount of bandwidth required to support it may become cost prohibitive. Cisco is in the process of addressing this concern by enhancing the CAC mechanisms provided by CUCM (Locations and RSVP) to support TelePresence. This functionality is scheduled for a future release of CUCM. As information about these enhancements becomes available, this document will be revised appropriately.



CHAPTER 9

Call Processing Deployment Models

Overview

For the current release of the Cisco TelePresence Solution and the Intra-Enterprise Deployment Model, a single Cisco Unified Communication Manager (CUCM) cluster is recommended to support all TelePresence devices within the enterprise. TelePresence meetings currently can only be scheduled across a single cluster by the Cisco TelePresence Manager (CTS-MAN) scheduling server because CTS-MAN only supports a single CUCM cluster. Although devices can register across multiple CUCM clusters, and ad hoc and permanent meetings can be established between clusters, this design is not currently recommended for customers deploying CTS-MAN. For customers not deploying CTS-MAN, this restriction is not applicable. Furthermore, a future release of CTS-MAN is planned to support multiple CUCM clusters, at which point this restriction will be removed.

In addition, in environments where TelePresence is deployed along with other generic Videoconferencing/Video Telephony devices on the same cluster, CUCM cannot instruct Videoconferencing/Video Telephony to use the recommended AF41 QoS marking and TelePresence to use the recommended CS4 QoS marking. The marking of audio and video traffic by CUCM is handled at the cluster level and not at the device level, because the marking of audio and video traffic is a cluster-wide (i.e., global) parameter and CUCM offers only a single parameter for video, which by default is set to AF41. For this reason it is recommended that TelePresence be placed on a separate cluster from all other Videoconferencing / Video Telephony applications. Finally, Cisco TelePresence requires CUCM release 5.1.1 or higher, with version 5.1.2 recommended to support the Auto Collaborate endpoint feature of TelePresence. Therefore, to summarize the guidance based upon the above three criteria, if a customer has a single existing cluster running version 5.1.1 or higher deployed for IP telephony and has no other Videoconferencing/Video Telephony devices, it is acceptable to integrate TelePresence devices onto that cluster. However, since the vast majority of deployments are not expected to meet these criteria, it is recommended that a separate CUCM cluster be deployed to support TelePresence and the guidance contained in this document is based upon that approach.

Dial-Plan Recommendations

For the current release of TelePresence, it is recommended that the Cisco Unified 7975G IP phones that serve as the user interface to the Cisco TelePresence system endpoints be marked to indicate that they should not be used for emergency services calls. A separate IP Phone registered to the production IP Telephony CUCM cluster should be deployed in the same room to provide access to emergency services.

To support functionality such as the ability to bridge audio participants into the TelePresence meeting via the audio add-in feature of the TelePresence System, the CUCM cluster which supports the TelePresence deployment may require additional components: either one or more voice gateways

connecting the TelePresence CUCM cluster to the customer's PBX or to the PSTN, and/or one or more Inter-Cluster Trunks (either H.323 or SIP) between the TelePresence CUCM cluster and the existing IP Telephony CUCM cluster(s).

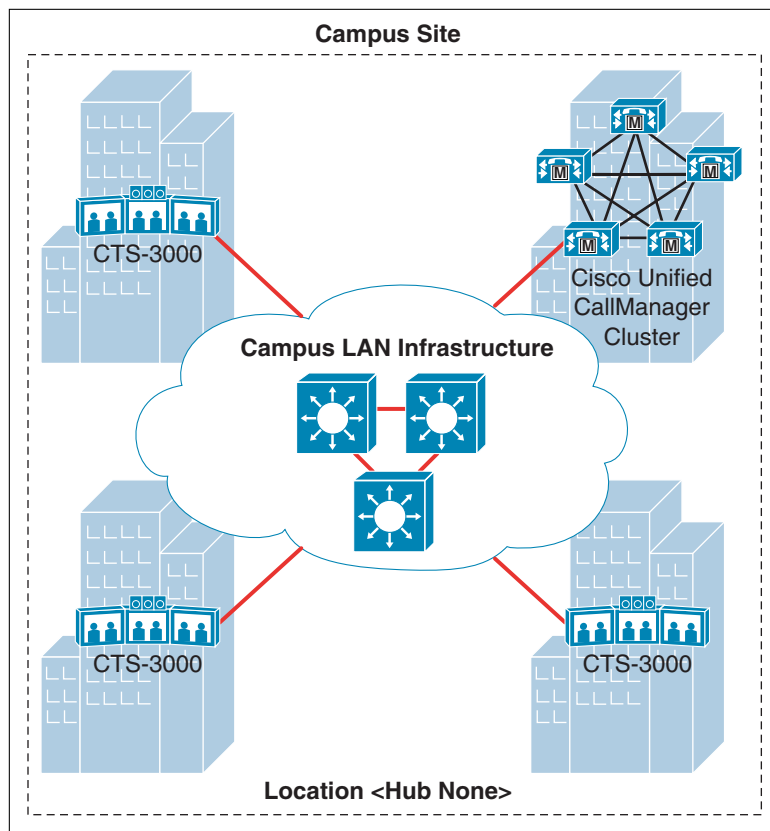
In either scenario, the TelePresence dial plan must be selected carefully and call routing set up appropriately to allow the TelePresence systems to reach and to be reached by other phones, audio conferencing bridges, and the PSTN. Therefore, the dial plan, Directory Numbers, Partitions, and Calling Search Spaces allocated to the TelePresence systems should be consistent with the rest of the enterprise to provide full support for current and future capabilities.

All current TelePresence deployments use either a single-site call processing model or a multi-site WAN with centralized call processing model. In both of these models, the CUCM cluster which supports the TelePresence devices resides at one location, such as a main campus. All communications with devices at remote locations takes place over the IP network infrastructure.

Single-Site Call Processing Model

The single-site call processing model applies to Cisco TelePresence deployments within a single campus and to deployments across MANs with LAN speed (i.e., Gigabit Ethernet) connectivity between sites. [Figure 9-1](#) shows an example of this deployment model.

Figure 9-1 Cisco TelePresence Single-Site Deployment



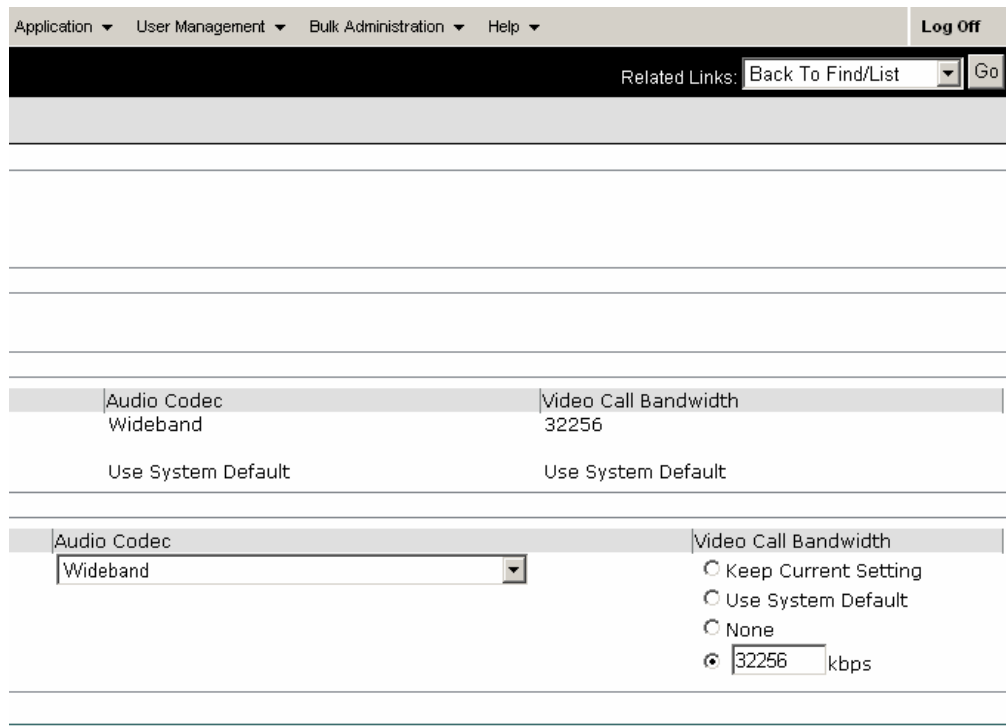
220218

Call Admission Control

In a single-site design, it is assumed that a high-speed LAN provides connectivity between all devices. CAC is typically not an issue, since the LAN can easily be scaled to provide sufficient bandwidth to simultaneously support all possible TelePresence meetings. TelePresence devices can be left within the default Hub_None location within the CUCM configuration, which provides no bandwidth restrictions on the total amount of video and audio traffic.

The region settings within the CUCM configuration are used to control the audio codec and the amount of video bandwidth used per call within a region and between regions. Since there are no other video devices in a standalone TelePresence deployment, all TelePresence devices can be placed in a single region. The region should be configured for AAC/Wideband audio (which as of release 5.1.1 or higher of CUCM permits up to 256 Kbps of audio per call) and a video bandwidth of at least 12500 Kbps (12.5 Mbps). As of release 5.1.1 or higher of CUCM, the maximum video bandwidth permitted is 32,256 Kbps. These settings are illustrated in Figure 9-2.

Figure 9-2 Recommended CUCM Region Settings for TelePresence

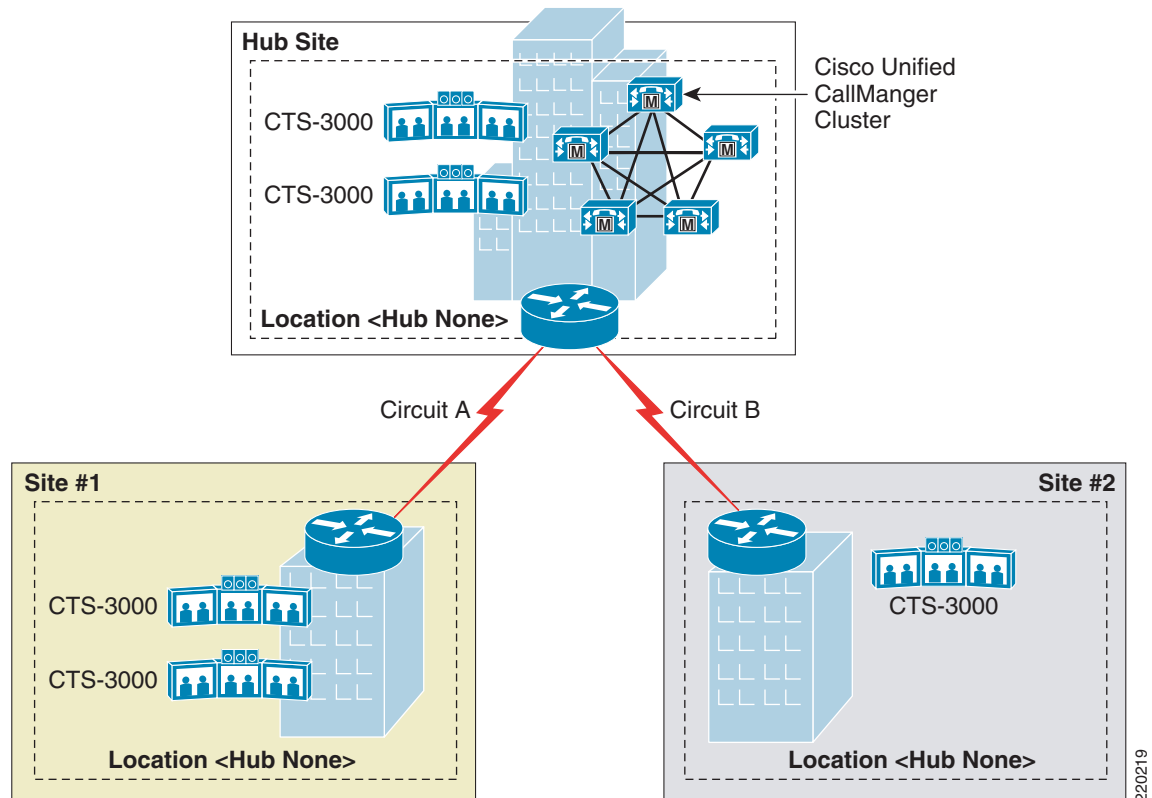


221699

Multi-Site WAN with Centralized Call Processing Model

In a multi-site WAN with centralized call processing model, a single CUCM cluster is deployed at a central site. This acts as the call processing agent for TelePresence devices both at the local and remote sites. Figure 9-3 shows an example of this deployment model over a hub-and-spoke network topology.

Figure 9-3 Cisco TelePresence Multi-Site Deployment

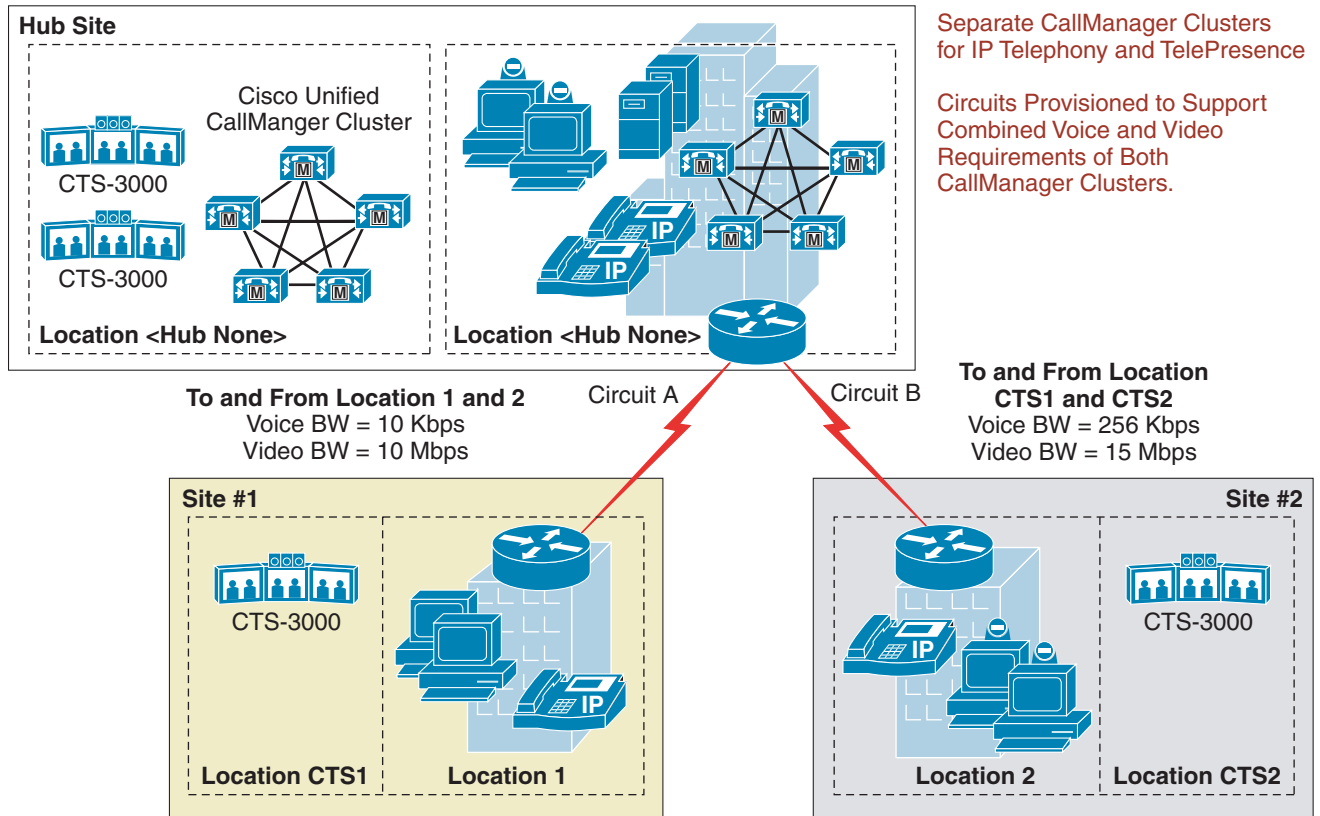


Call Admission Control

For current TelePresence deployments it is recommended that sufficient WAN bandwidth be provisioned to support all possible simultaneous meetings within the network. Refer to [Chapter 8, “Capacity Planning and Call Admission Control”](#) for details regarding the use of manual capacity planning to guarantee 100% call completion. For this design, all TelePresence devices can be left in the default Hub_None location which provides no bandwidth restrictions on the total amount of video and audio traffic (as shown above). Alternatively, TelePresence devices at each remote site can be assigned to a different location and the video and audio bandwidth between locations set to unlimited. It should be noted that when deploying multiple CTMS devices in a distributed multipoint TelePresence deployment, TelePresence devices need to be assigned to separate regions/locations in order for CTS-MAN to correctly choose the appropriate CTMS for the multipoint meeting. Further details are discussed in [Chapter 10, “Cisco TelePresence Multipoint Solution Essentials,”](#) [Chapter 11, “Cisco Multipoint Technology and Design Details,”](#) and [Chapter 12, “Cisco TelePresence Multipoint Solution Circuit and Platform Recommendations.”](#)

When implementing Cisco TelePresence alongside an existing CUCM deployment dedicated for IP telephony, the WAN circuits must be provisioned with sufficient bandwidth to take into account the CAC requirements of both CUCM clusters. An example of this is shown in Figure 9-4.

Figure 9-4 Separate Cisco Unified CUCM Design Example



As can be seen in Figure 9-4, separate CUCM clusters are deployed for TelePresence and for IP telephony (both dashed boxes). Each CUCM configuration has a different location configured for each remote site with a certain amount of bandwidth configured between each location for audio and video. In this scenario, the WAN circuits must be provisioned to accommodate the aggregate bandwidth pools configured in both CUCM clusters, since they operate independently of each other. Otherwise, the potential exists for oversubscribing the circuits and degrading the quality of voice, desktop video, and TelePresence meetings.

It should also be noted that the Survivable Remote Site Telephony (SRST) feature of Cisco router platforms do not currently support Cisco TelePresence system devices. Therefore in a multi-site WAN with a centralized call processing TelePresence design, SRST cannot be used to provide redundancy if the connection to the TelePresence CUCM cluster fails. However in the design shown in Figure 9-4, where a separate CUCM cluster is deployed for IP telephony devices, SRST works well for the IP phones and other devices which are supported.



CHAPTER 10

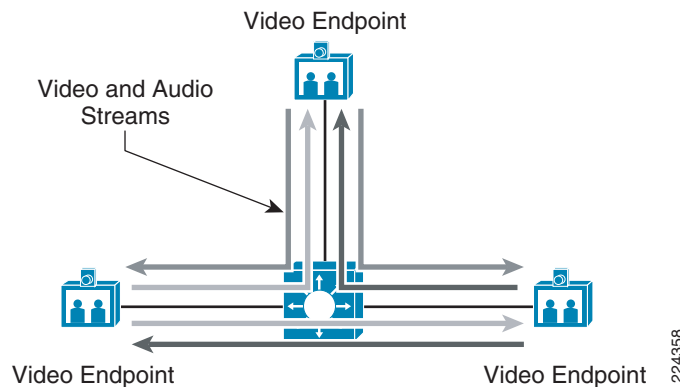
Cisco TelePresence Multipoint Solution Essentials

Overview of Multipoint Conference Technologies

This section briefly discusses some of the methods and technologies used to currently provide multipoint conferencing.

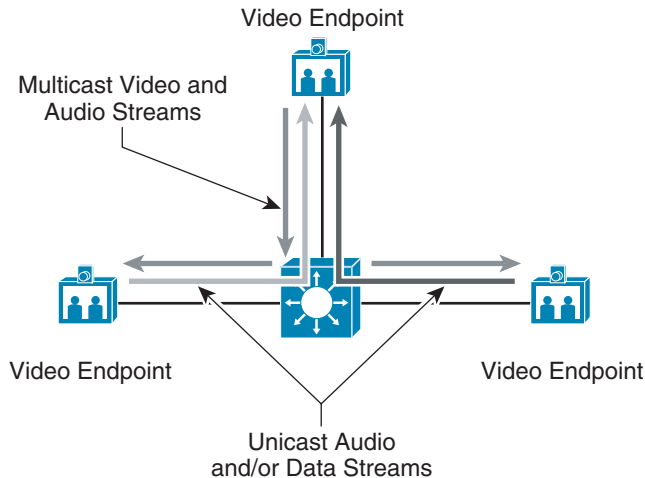
Figure 10-1 shows the simplest form of multipoint conferencing in which each video endpoint (also known as a video terminal when referring to H.323 based-systems) simply sends unicast video and audio streams to every other video endpoint in the multipoint call.

Figure 10-1 Multipoint Conferencing Using Point-to-Point Unicast Streams



This form of multipoint video conferencing has limited scalability. Each endpoint has to mix the audio and video from every other endpoint in order to determine which video stream to display at any given time. Alternatively, the endpoints may display video from each of the other endpoints in a small section of the display, for example four or eight boxes of video. Audio must still be mixed at the video endpoint and played out on the speaker(s). The addition of each new video endpoint results in every other endpoint having to both send and receive a new video and audio stream. For a conference of N video endpoints, each endpoint has to both send and receive $N-1$ unicast video and audio streams. Bandwidth utilization is somewhat inefficient for this reason.

Figure 10-2 shows a second form of multipoint conferencing in which an endpoint multicasts its video and audio stream to every other endpoint in the multipoint call.

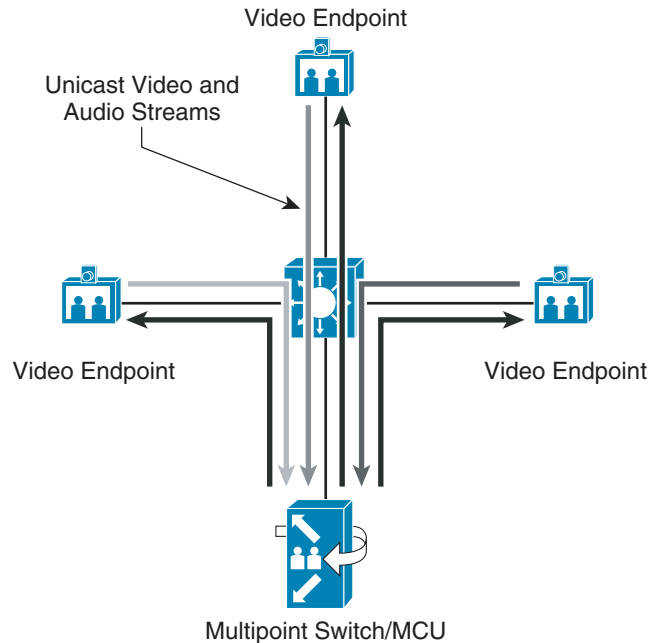
Figure 10-2 One-Way Multipoint Conferencing Using a Multicast Stream

224359

This form of multipoint conferencing requires multicast friendly media, such as a LAN/WAN infrastructure which is IP multicast enabled. It is not possible to implement this over point-to-point media such as ISDN. Typically this form of video conferencing is used for one-way multipoint conferences, with a single endpoint functioning as the “presenter” and the remaining endpoints functioning as the “audience.” The presenter may have a camera, but no video display. Likewise the audience may have video displays, but no cameras. A unicast data path or audio stream (often through a voice conferencing bridge) is sometimes utilized to allow the audience to ask questions and provide feedback to the presenter.

Full two-way video conferencing using IP multicasting is not extensively utilized. Each video endpoint would need to multicast to a separate IP address to achieve this. Although less inefficient in terms of network utilization than point-to-point unicast streams, it suffers from the same scalability issue in that each video endpoint has to mix the audio and video streams from every other endpoint in order to determine which video stream to display at any given time.

Figure 10-3 shows a third form of multipoint conferencing in which each video endpoint sends unicast audio and video streams to a centralized multipoint switch or multipoint control unit (MCU).

Figure 10-3 Multipoint Conferencing Using a Multipoint Switch / MCU

The multipoint switch/MCU mixes the audio and video streams from each video endpoint and transmits a single audio and video stream back to each endpoint. It should be noted that if multiple displays or speakers exist at an endpoint, multiple video and audio streams may be sent by the multipoint switch/MCU to that endpoint.

If audio or video transcoding or transrating are required, a more traditional MCU function can be deployed which decodes and re-encodes the audio or video streams to each endpoint. This can lead to somewhat higher latency (end-to-end delay) of the audio and video. When transcoding or transrating are not required—as with the Cisco TelePresence solution—a multipoint switch which simply mixes and switches audio and video streams to each endpoint can be deployed. This leads to much lower latency. Since the mixing and potentially transcoding or transrating functionality has been offloaded to the multipoint switch/MCU, scalability can be greatly enhanced by deploying custom-built platforms for this functionality. Enhancements to this technology also make use of voice activity switching, in which voice activity from the audio streams of the endpoints is used to signal which endpoints should be sending video to the multipoint switch/MCU. In this way, only the endpoints which are actively speaking send video to the multipoint switch/MCU, which then transmits the video and audio to the other endpoints.

The Cisco TelePresence suite of virtual meeting solutions utilizes this third method of providing multipoint virtual meetings via the Cisco TelePresence Multipoint Switch (CTMS).

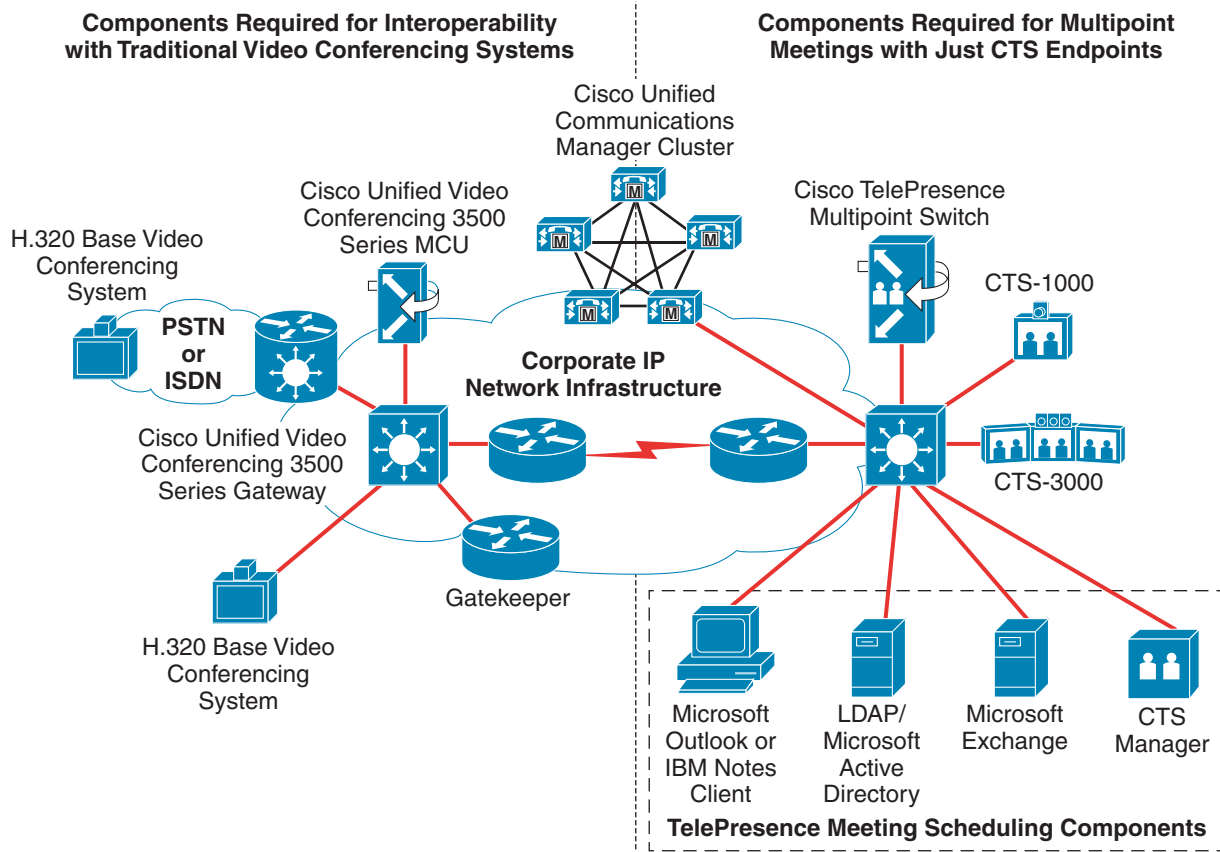
Components of the Cisco TelePresence Multipoint Solution

This section presents the hardware and software components required to deploy multipoint TelePresence within the network infrastructure.

Overview

Figure 10-4 shows the components required for a Cisco multipoint TelePresence virtual meeting.

Figure 10-4 Components of a Cisco TelePresence Multipoint Meeting



Depending upon whether the TelePresence virtual meeting includes only Cisco TelePresence Systems (CTS) endpoints or also includes traditional video conferencing endpoints, the required components differ. The following sections list the components for each type of meeting.

Multipoint Virtual Meetings Which Include Only CTS Endpoints

The right side of Figure 10-4 shows the components required for a TelePresence meeting which consists of CTS endpoints only. The components are:

- Cisco TelePresence Multipoint Switch
- Three or more CTS endpoints—These CTS endpoints can consist of any combination of CTS-500, CTS-1000, CTS-3000, or CTS-3200 units.
- A Cisco Unified Communications Manager (CUCM) Cluster—This provides call signaling and control for the CTS endpoints.
- An IP infrastructure for transport of the call signaling, voice, and video media.
- Optional Cisco TelePresence meeting scheduling components—These consist of the CTS Manager (CTM), Microsoft Exchange or IBM Domino Server, LDAP/Microsoft Active Directory Server, and PCs running Microsoft Outlook or IBM Notes client.

Multipoint Virtual Meetings Which Also Include Traditional Video Conferencing Systems

If traditional video conferencing systems (non-Cisco TelePresence endpoints) are required to attend the virtual meeting, the components on the left side of [Figure 10-4](#) may also be required, in addition to the components listed on the right side. They are:

- A Cisco Unified Video Conferencing 3515 MCU or a Cisco Unified Video Conferencing 3545 System functioning as an MCU.
- One or more traditional video conferencing systems (LAN-, PSTN-, or ISDN-based).
- Optional Cisco gateway for interoperability between H.323 (LAN-based) and H.320 (ISDN-based) or H.324 (POTS based) video conferencing systems.
- Optional Cisco IOS gatekeeper (ISR router) for E.164 address resolution, call routing, and call admission control of H.323 based systems.

This document only discusses multipoint TelePresence meetings which consist of Cisco TelePresence endpoints. Future versions may include interoperability with traditional video conferencing systems.

Cisco TelePresence Multipoint Switch Overview

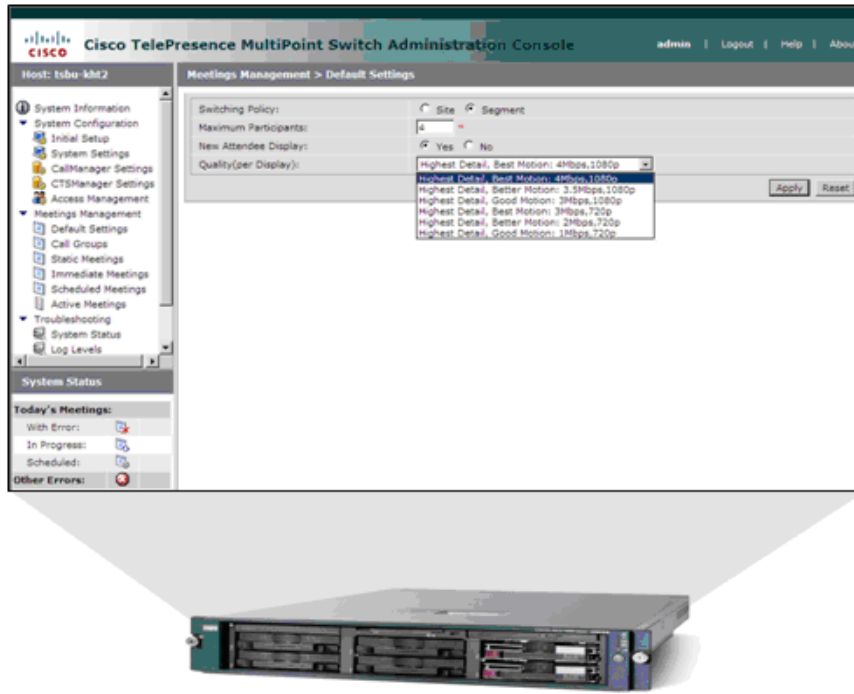
The Cisco TelePresence Multipoint Switch (CTMS) is a purpose built multipoint appliance developed by Cisco to directly address multipoint requirements for Cisco TelePresence. The patent pending software architecture of the CTMS provides extremely low-latency video and audio switching, adding less than 10ms of delay to any multipoint meeting. In combination with the CTS Manager, CTMS provides a scalable multipoint solution for a Cisco TelePresence network of any size.

CTMS provides a very scalable architecture supporting up to 48 simultaneous table segments—a table segment being defined as a display and camera on any CTS system. CTS-3000 and CTS-3200 systems consist of three table segments, while CTS-1000 and CTS-500 systems consist of a single table segment. Each CTMS provides support for up to 16 CTS-3000s or CTS-3200s, 48 CTS-1000s or CTS-500s, or any combination up to 48 table segments in a single meeting or any number of meetings.¹

The high-performance, server-based architecture and Linux-based Cisco Voice Operating system of the CTMS provides a familiar and reliable platform. CTMS provides system management via Secure Shell (SSH), Hyper-Text Transfer Protocol over Secure Sockets Layer (HTTPS) and Cisco Discovery Protocol (CDP). From an administrator's perspective, CTMS is managed using tools and methodologies that are consistent with those used for Cisco Unified Communications Manager, CTS Manager, and CTS rooms.

1. CTMS version 1.0 supports a maximum of 36 table segments. CTMS version 1.1 extends this to a maximum of 48 table segments.

Figure 10-5 Cisco TelePresence Multipoint Switch



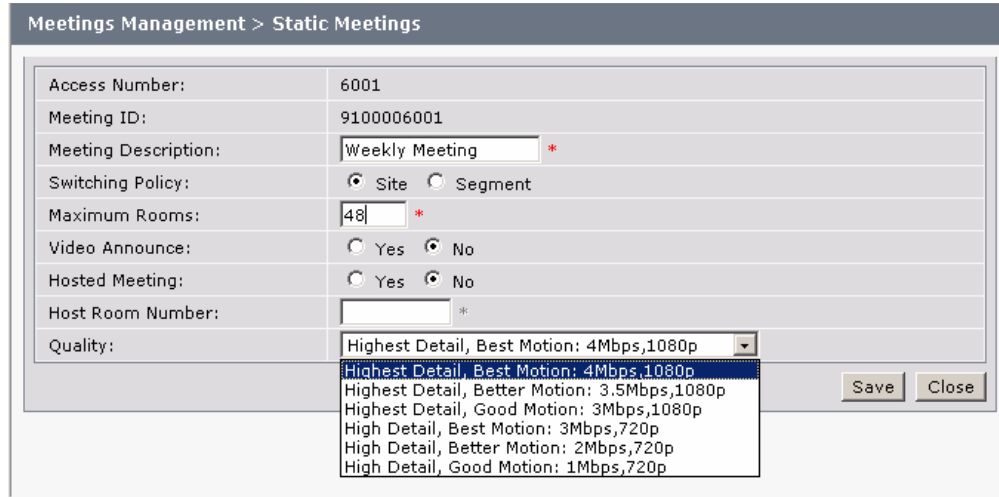
Meeting Types

The CTMS supports multiple meeting types, providing users with a number of scheduled and non-scheduled meetings to choose from based on their meeting requirements. CTMS supports scheduled meetings, non-scheduled (static or ad hoc) meetings, or a combination of both meeting types in a single deployment environment. CTS Manager is required for scheduled meetings. It provides CTMS resource management and integration with Microsoft Exchange or IBM Domino. All supported meeting types are defined below.

Static Meetings

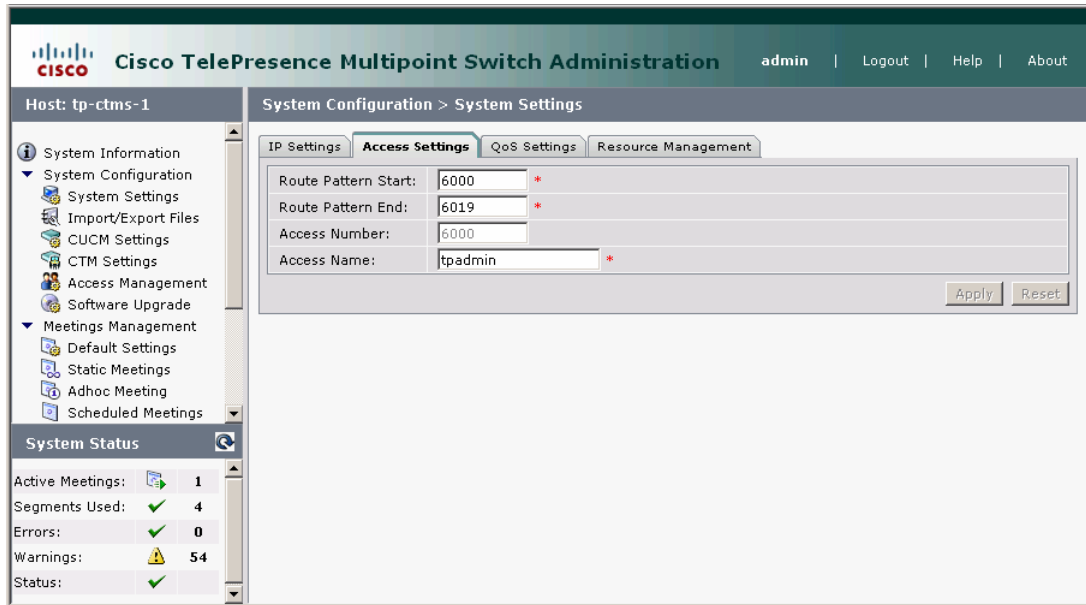
Static meetings are non-scheduled meetings configured on the CTMS through the administration GUI. An example is shown in [Figure 10-6](#).

Figure 10-6 Static Meeting Configuration



A meeting scheduler or administrator, who sets up the static meeting, manually assigns a meeting access number that is used to access the meeting. Meeting access numbers must be chosen from the range of numbers between the Route Pattern Start and Route Pattern End configured on the Access Settings page of the CTMS, as shown in Figure 10-7.

Figure 10-7 Dial Pattern Range Configured on the CTMS



The start and stop of the route pattern should correspond to the range of numbers configured within CUCM for the SIP trunk corresponding to the CTMS. Static meetings are always available to any CTS room. They are accessed by manually dialing the static meeting telephone number or using a speed dial entry on the Cisco 7975G IP Phone associated with the CTS device within the TelePresence room. Static meetings are the equivalent of a meet-me meeting in the voice world.

A meeting scheduler or administrator may add additional CTS rooms to the meeting at any time using the CTMS administrative GUI. The meeting scheduler or administrator also determines the meeting features for the static meeting, each of which is discussed in the following sections.

Static meetings may be hosted or non-hosted. Hosted static meetings require the meeting scheduler or administrator to pre-assign a host room to be present in order to start the meeting. The telephone number of the host room is configured within the static meeting. All meeting rooms which dial-in before the host room are placed on hold until the host room joins the meeting. All rooms are disconnected once the host room leaves the meeting.

Ad Hoc Meetings

Ad hoc meetings are non-scheduled, administrator-initiated, dial-out meetings. A meeting scheduler or administrator initiates the meeting through the CTMS administrative GUI by listing the telephone number of the rooms which will participate in the multipoint meeting, as shown in [Figure 10-8](#).

Figure 10-8 Ad Hoc Meeting Configuration

The screenshot displays the Cisco TelePresence Multipoint Switch Administration interface. The top navigation bar includes the Cisco logo, the title 'Cisco TelePresence Multipoint Switch Administration', and user options: 'admin | Logout | Help | About'. The left sidebar shows a tree view with categories: System Information, System Configuration (System Settings, Import/Export Files, CUCM Settings, CTM Settings, Access Management, Software Upgrade), Meetings Management (Default Settings, Static Meetings, Adhoc Meeting, Scheduled Meetings), and System Status. The main content area is titled 'Meetings Management > Adhoc Meeting' and contains a 'New Adhoc Meeting' form. The form has two tabs: 'New Adhoc Meeting' (active) and 'Meeting Templates'. Fields include: Meeting Template (dropdown), Rooms (text input with a list of '1010', '1011', '1011' and a note to 'Enter multiple numbers separated by Carriage Return (the ENTER Key)'); Meeting Description (text input with 'Adhoc Meeting'); Switching Policy (radio buttons for Site and Segment, with Segment selected); Video Announce (radio buttons for Yes and No, with Yes selected); and Quality (dropdown menu with an open list of options: 'Highest Detail, Best Motion: 4Mbps,1080p', 'Highest Detail, Better Motion: 3.5Mbps,1080p', 'Highest Detail, Good Motion: 3Mbps,1080p', 'High Detail, Best Motion: 3Mbps,720p', 'High Detail, Better Motion: 2Mbps,720p', and 'High Detail, Good Motion: 1Mbps,720p'). 'Apply' and 'Reset' buttons are at the bottom right. A vertical ID '224964' is on the right edge.

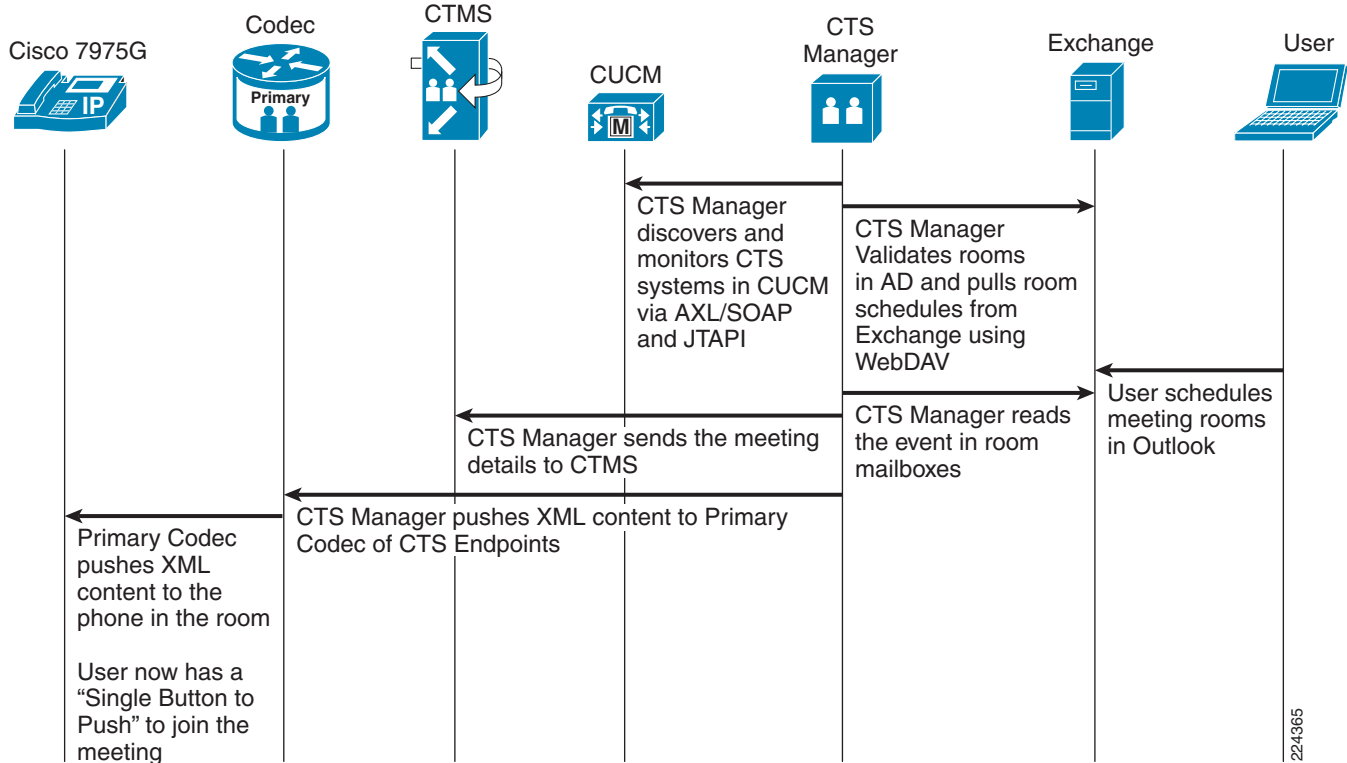
The CTMS dials out to each of the rooms, thereby requiring no end user interaction. Ad hoc meetings cannot be accessed by dialing in from a CTS system; rooms must be added by the meeting scheduler or administrator using the CTMS administrative GUI.

Scheduled Meetings

Multipoint TelePresence meetings are scheduled by end users using Microsoft Exchange or IBM Domino clients in the same manner that point-to-point meetings are scheduled. Scheduled meetings require no CTMS administrator interaction. CTS Manager is a required component for scheduled meetings. It provides the interface between Microsoft Exchange or Lotus Domino and the CTMS, allowing the appropriate resources on the CTMS to be reserved for the multipoint meeting.

[Figure 10-9](#) shows the interaction between the CTS Manager, Microsoft Exchange, CTMS, and the CTS endpoints when scheduling a multipoint TelePresence meeting.

Figure 10-9 Information Flow When Scheduling a Multipoint Meeting



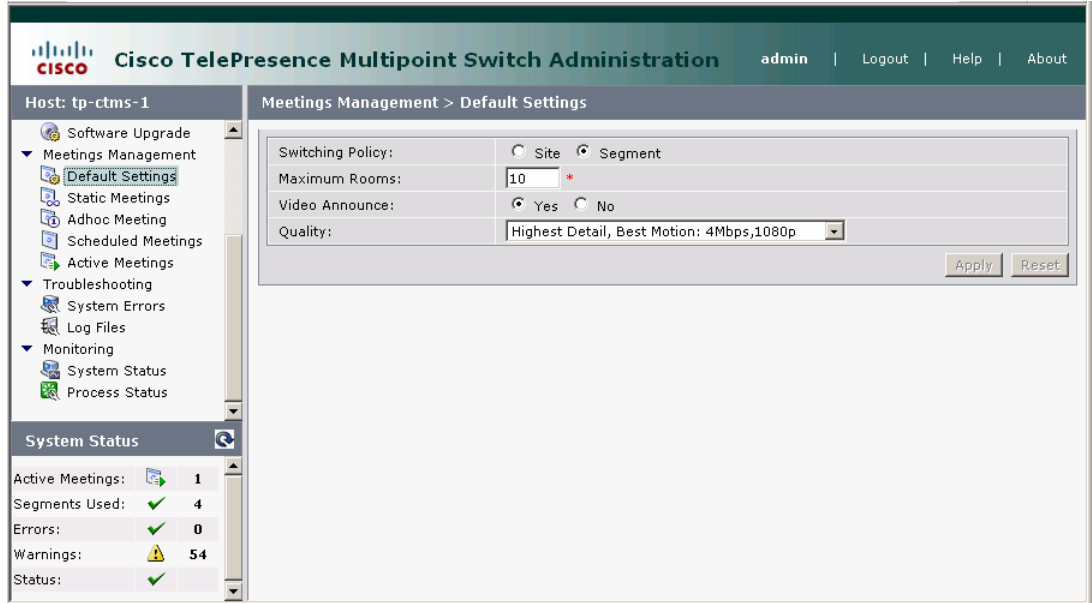
CTS Manager uses the WebDav protocol to communicate with Exchange in order to update scheduled meetings. CTS Manager also uses AXL/SOAP and JTAPI to discover and monitor the CTS endpoints defined within CUCM. When an end user schedules a meeting within Outlook, CTS Manager reads the meeting events scheduled within the mailboxes defined for each CTS room within Exchange. CTS Manager sends the meeting details to CTMS via XML in order to allocate the necessary resources from the Schedulable Segments resource pool. The primary codec of the CTS endpoint then updates the schedule on the associated IP 7975G phone via XML as well. Finally, the CTS Manager confirms the meeting reservation to the end user, once the CTMS and CTS endpoint resources have been confirmed. It does this via the mailbox configured for the CTS Manager device within Exchange.

Scheduled meetings can be accessed by using the one-button-to-push meeting access feature. Additional rooms can be added to active scheduled meetings at any time during the meeting by the meeting scheduler or administrator, using the CTMS administrative GUI. It should be noted that additional rooms added by a meeting scheduler or administrator count against the Ad hoc Segments resource pool and not the Schedulable Segments resource pool. [Multipoint Resources](#) discusses these resource pools.

CTMS Meeting Features

The CTMS supports a number of meeting features that provide flexibility and security for multipoint meetings. Each multipoint meeting is configured with a default meeting policy that may be modified by the system administrator before or during the multipoint meeting. [Figure 10-10](#) shows an example of the default meeting settings for a CTMS.

Figure 10-10 CTMS Default Meeting Settings



When individual meetings are configured, the default features can be overridden. The following sections discuss these features in detail.

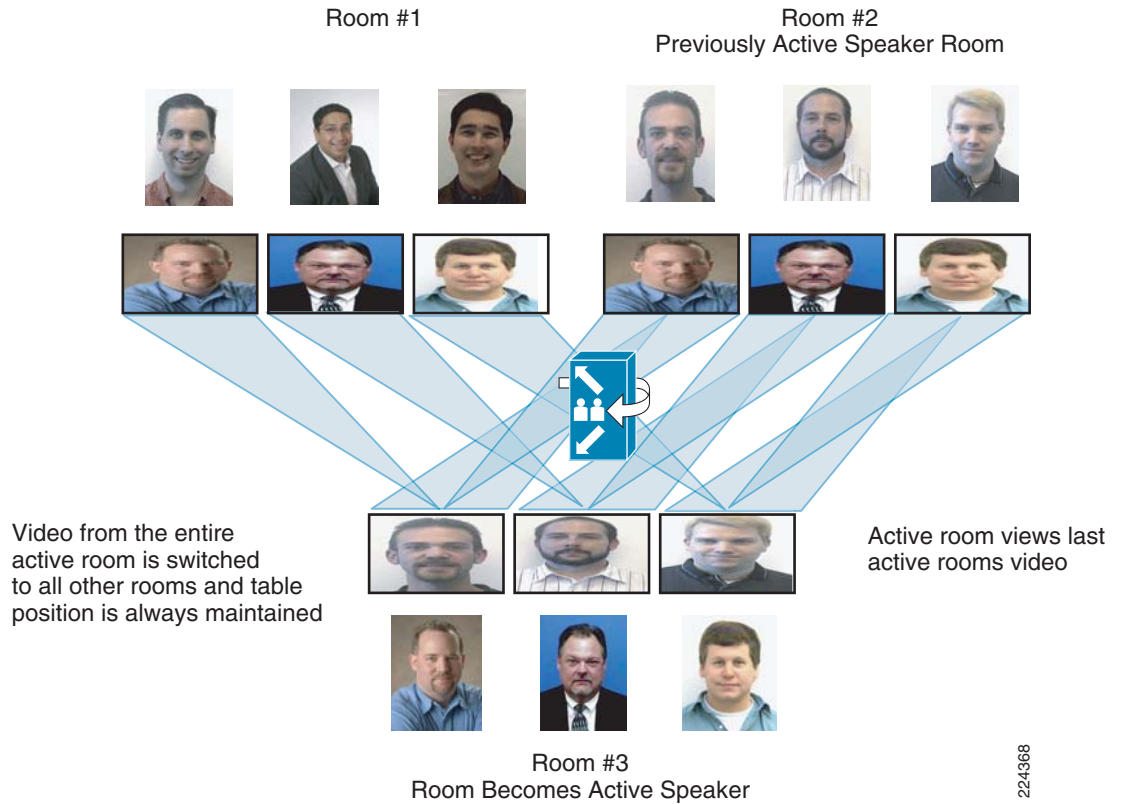
Switching Policy

The switching policy has two choices, room switching and speaker switching. Room switching applies to multipoint calls which include CTS endpoints with multiple cameras and displays, such as the CTS-3000 and the CTS-3200. Speaker switching applies to all CTS endpoints.

Room Switching

Room switching switches the video from all table segments of a particular room to all other rooms in a multipoint meeting. For CTS endpoints with multiple screens, if the active speaker (loudest speaker for approximately two seconds) changes, all table segments in the new active speaker's room are displayed in all other rooms at the same time, replacing the previously active speaker's room. The position of each of the screens—left, center, or right—on the active speaker's room is maintained on each of the other rooms. Figure 10-11 shows an example of room switching with CTS-3000s.

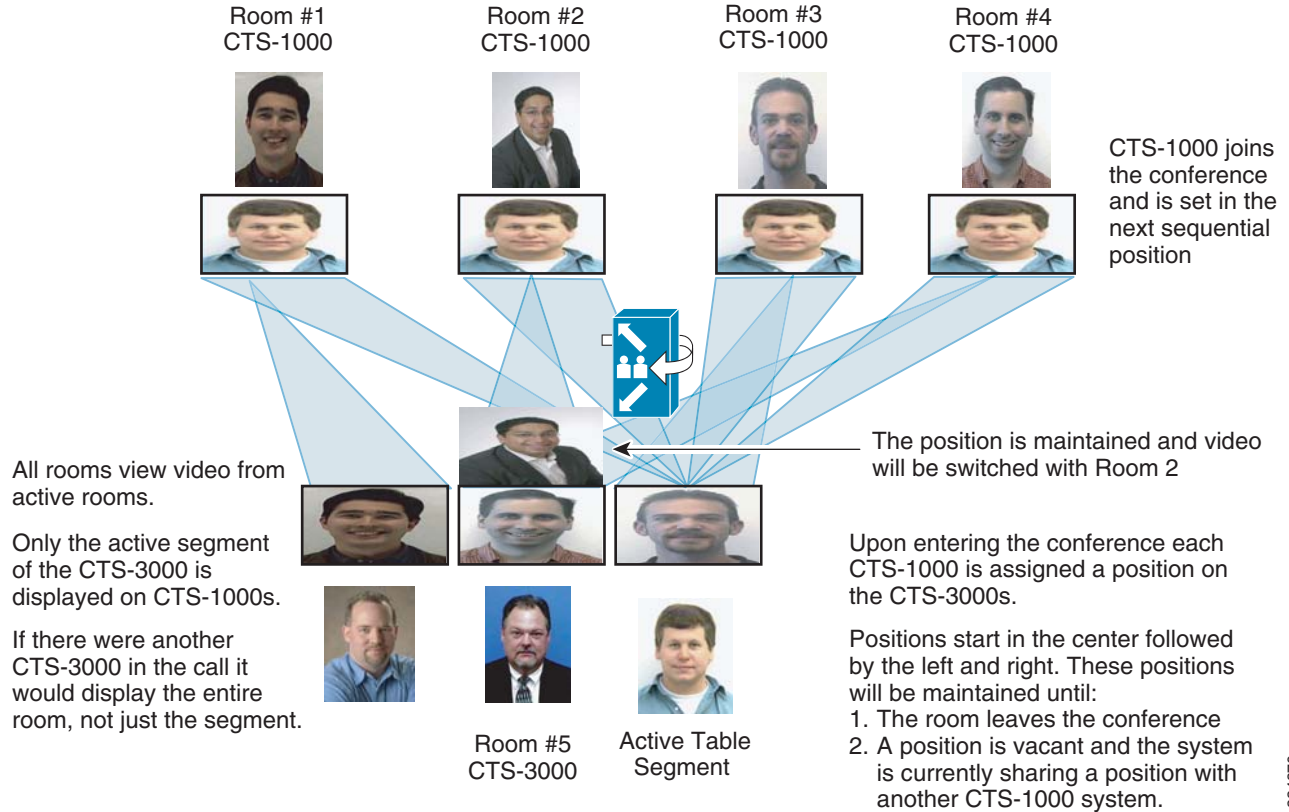
Figure 10-11 Room Switching with CTS-3000s



As can be seen, Room 3 has become the active speaker. Video is displayed from the cameras in Room 3 to both of the other rooms. Since Room 2 was the previously active speaker room, video from Room 2 is still displayed on the screens of Room 3.

Room switching can also be applied to multipoint conferences which include both CTS-3000s and CTS-1000s. This is shown in [Figure 10-12](#).

Figure 10-12 Room Switching with CTS-1000s and CTS-3000s



224370

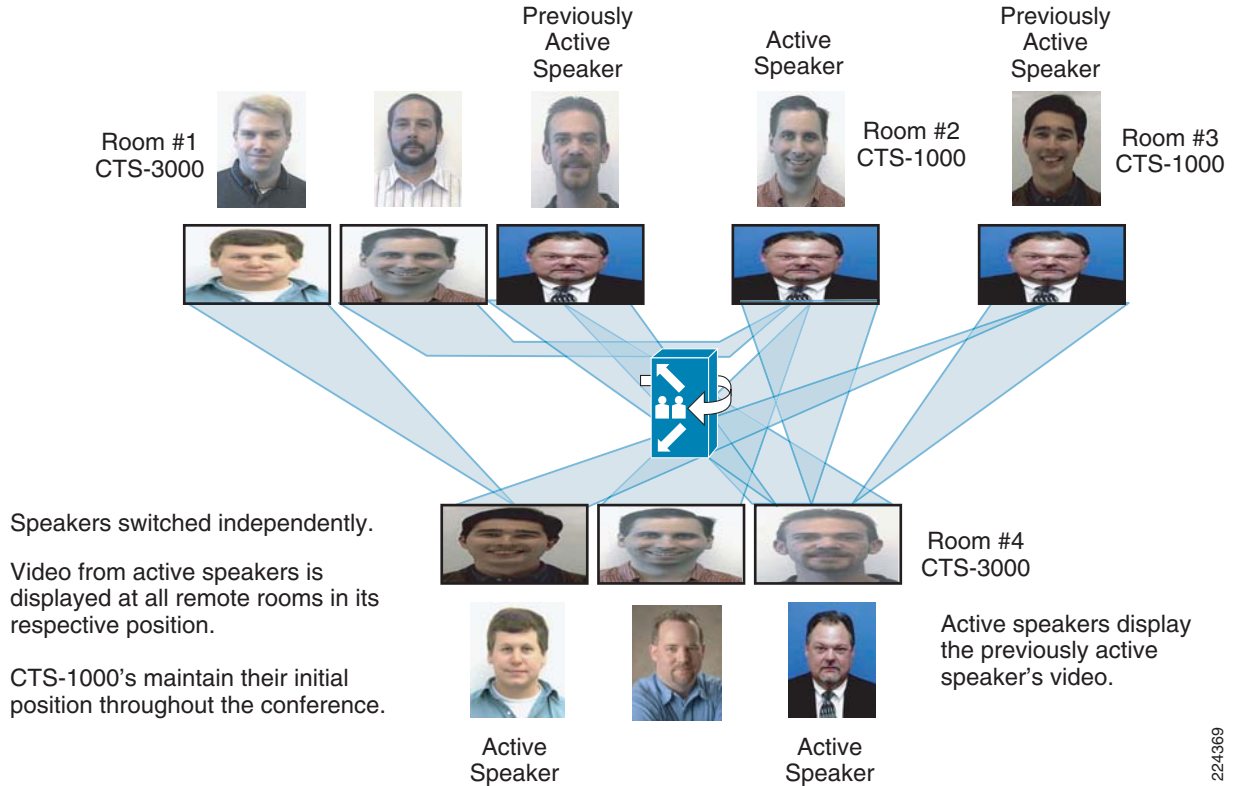
In this case, when a particular table segment at a CTS-3000 room becomes the active speaker, only that table segment is displayed on all CTS-1000 rooms. However, if there were other CTS-3000s in the call, all three table segments from the active speaker's room would be displayed on the other CTS-3000 rooms.

Also notice in [Figure 10-12](#) that as each CTS-1000 enters the multipoint meeting, it is assigned a display position on the CTS-3000s within the meeting. Positions start with the center, followed by the left, and then right. The position of a particular CTS-1000 is maintained until either the room leaves the conference, or a different position on the CTS-3000s is vacant and the system is currently sharing the current CTS-1000 position with another CTS-1000. CTS-1000 rooms which share positions on CTS-3000 rooms are switched by the CTMS based upon voice activity from the particular room.

Speaker Switching

Speaker switching allows each table segment to be switched independently, so that at any given time a room may be viewing three different rooms. The new active speaker's table segment is displayed in all other rooms in its proper position. [Figure 10-13](#) shows an example of speaker switching with CTS-1000s and CTS-3000s.

Figure 10-13 Speaker Switching with CTS-1000s and CTS-3000s



224369

As each CTS-1000 enters the multipoint meeting, it is again assigned a display position on the CTS-3000s within the meeting. The position of a particular CTS-1000 is maintained until the room leaves the conference.

Room Switching versus Speaker switching has implications regarding the bursts generated on the network when a new room or table segment becomes the active speaker. Generally speaker switching is less bursty than room switching. This is discussed in more detail in [Estimating Burst Sizes within Multipoint TelePresence Calls](#).

Maximum Number of Rooms

The Maximum Rooms setting allows the administrator to limit the number of rooms allowed to join the meeting. The maximum number of rooms supported in a single call, and by the CTMS overall, is currently 48 rooms as of CTMS version 1.1. Prior to version 1.1 the maximum number of rooms supported was 36. It should be noted that Maximum Rooms actually refers to the maximum number of segments supported, where a CTS-3200 or CTS-3000 counts as three segments and a CTS-1000 or CTS-500 counts as one segment. The Maximum Rooms setting works hand-in-hand with CTMS resource management as discussed in [Multipoint Resources](#).

Video Announce

The Video Announce feature causes a Cisco TelePresence room which joins the conference to be displayed to all other rooms for approximately two seconds. This prevents a muted room from joining without being noticed.

The Video Announce feature also has implications regarding the bursts generated on the network when a new site joins the conference. This is discussed in more detail in [Estimating Burst Sizes within Multipoint TelePresence Calls](#).

Lock Meeting

Administrators can lock a meeting through the administrative GUI after a meeting is in progress. This prevents any additional rooms from joining the meeting. Additional TelePresence systems can be added to a meeting that is locked, but only by the meeting scheduler or administrator through the administrative GUI.

Quality

The quality setting allows the configuration of 1080p or 720p on a per call basis. It is important that meetings are configured to support the lowest resolution Cisco TelePresence room participating in the multipoint meeting because the CTMS does not perform any transcoding or transrating.

Cisco TelePresence systems have the ability to negotiate down from 1080p to 720p, allowing systems configured for 1080p to join a meeting configured for 720p. However, Cisco TelePresence systems configured for 720p cannot negotiate up to 1080p and therefore will not connect to a meeting configured for 1080p.

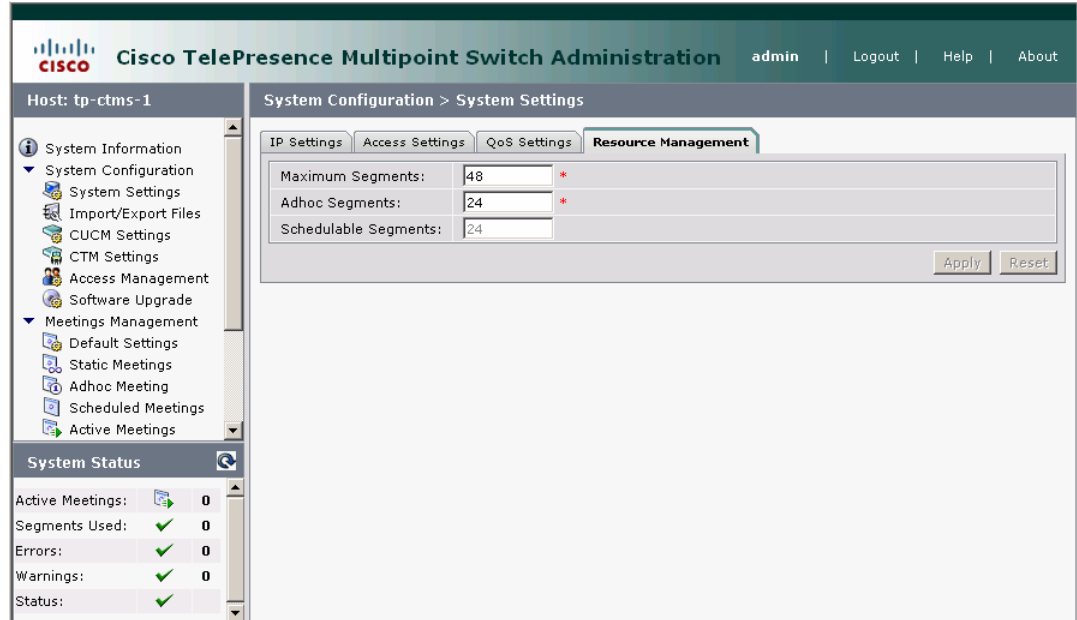
VIP Mode

VIP Mode is a switching feature which is useful for meetings in which one site or segment of a site is more important than the rest of the sites. VIP Mode can be configured using either a “hard-lock” or “soft-lock.” When using a “hard-lock,” the video from the VIP site or segment is always displayed at the other sites, regardless of who is the active speaker. When using a “soft-lock,” the video from the VIP site or segment is temporarily switched when a speaker at another site talks, but is automatically switched back after the speaker stops talking, without the VIP site having to talk. This ensures the focus of the meeting is always on the VIP site.

Multipoint Resources

Multipoint resources are configured on the CTMS based on table segments. As mentioned above, each CTMS supports a maximum of 48 segments that can be allocated for scheduled or ad hoc meetings. Resources are defined in the CTMS administrative GUI as maximum segments, ad hoc segments, and schedulable segments. [Figure 10-14](#) shows the multipoint resource configuration options.

Figure 10-14 CTMS Resource Management Configuration



- Maximum segments—Total number of segments supported on a CTMS. This field may be used to limit the number of segments supported on the CTMS (48 maximum).
- Ad hoc segments—Number of segments available for non-scheduled meetings. Ad hoc segments are also used for any non-schedule CTS endpoint added to a scheduled multipoint meeting through the CTMS administrative GUI.
- Schedulable segments—Number of segments available for scheduled meetings. The number of schedulable segments is used by CTS Manager to track available segments for scheduled multipoint meetings. This number is automatically calculated by the system using maximum and ad hoc segments.

It is possible to define multiple static meetings, the total of which exceeds the Ad hoc Segment resource pool. Static meetings are basically permanent meetings that can be utilized at any time without the administrator being notified. Because ad hoc meetings and static meetings are deducted from the same resource pool, it is possible to have insufficient table segments available for an ad hoc meeting due to ongoing static meetings and vice versa. Further, as mentioned earlier, when a meeting scheduler or administrator adds another TelePresence room to a scheduled meeting, the resources for that room are deducted from the Ad hoc Segments resource pool. It is therefore possible to have insufficient resources to add the room to the scheduled meeting due to ongoing static or ad hoc meetings.

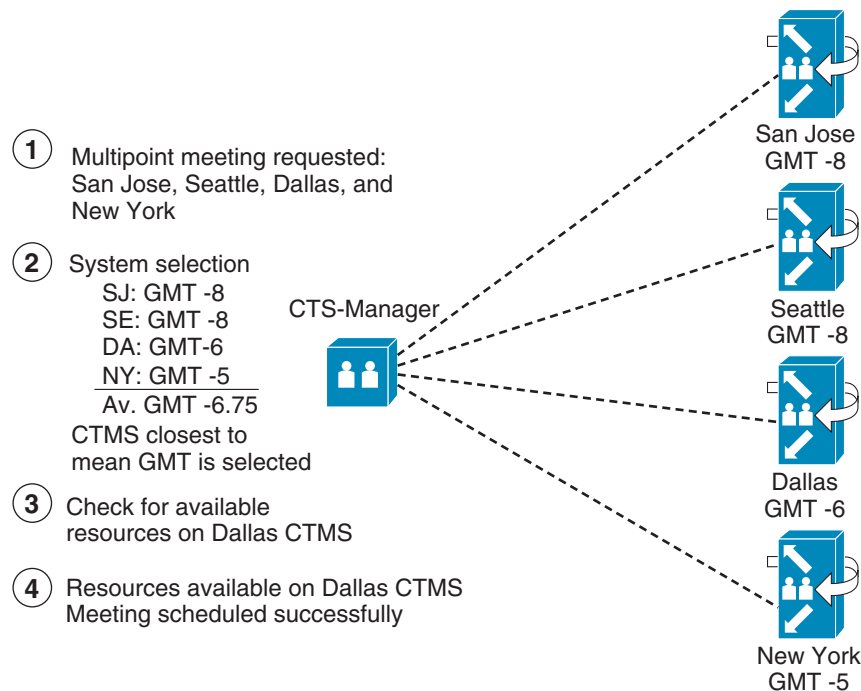
Because of these issues, it is recommended that in distributed CTMS deployments which support combined scheduled and non-scheduled meetings, separate CTMS units be deployed for static meetings. One set of CTMS units should be allocated for scheduled and ad hoc meetings and the second set for static meetings. This is discussed further in [Deployment Models](#).

Geographical Resource Management

When multiple CTMS devices are available in a multipoint deployment, it is important that multipoint meetings are hosted on the appropriate CTMS based on proximity to CTS endpoints. This helps regionalize multipoint meetings, conserve bandwidth between regions, and ensure the lowest latency for each meeting. Managing large, distributed multipoint deployments without geographical resource management is an unmanageable task.

Geographical resource management requires CTS Manager (CTS-MAN) and is only available for scheduled meetings. Currently geographical selection is based on time zones. For each scheduled multipoint meeting, CTS-MAN selects the CTMS whose GMT timezone is closest to the calculated mean GMT, checks for available resources on the selected CTMS, and schedules the meeting. If there are no resources available on the selected CTMS, CTS-MAN works its way down the list of CTMS devices until available resources are found. [Figure 10-15](#) illustrates the CTMS selection process.

Figure 10-15 CTMS Selection Example



Note: If no resources are available in Dallas the next closest CTMS is selected (San Jose GMT -8)

22/4584

The time zone of each CTS endpoint is determined by the Date/Time Group to which it is assigned within the CUCM configuration. However, CUCM devices are not directly assigned to a specific Date/Time Group. Instead, separate Device Pools must be created for each individual time zone. The CTS endpoints are then assigned to the Device Pool which holds the Date/Time Group corresponding to the time zone within which the endpoint is located. For a distributed multipoint design to function correctly, it is therefore necessary to have CTS units assigned to multiple Device Pools which reflect their correct geographic location. Note that this is not necessarily a requirement of a centralized multipoint design, since all CTS endpoints utilize the single CTMS regardless of time zone in such a deployment.

The time zone of the CTMS is configured directly on the CTMS under System Configuration-->CTM Settings, as shown in [Figure 10-16](#). Note that the SIP Trunk configuration for the CTMS within CUCM also has a Device Pool setting. The time zone configured within the CTMS should match the time zone configured within the CUCM Device Pool.

Figure 10-16 CTMS Location Configuration

The screenshot displays the Cisco TelePresence Multipoint Switch Administration interface. The main content area is titled 'System Configuration > CTM Settings'. The configuration fields are as follows:

Description:	CTMS *
Latitude:	35 47 N *
Longitude:	78 39 W *
Time Zone:	US/Eastern -- Eastern Standard Time *
Daylight Savings Time:	<input checked="" type="radio"/> Yes <input type="radio"/> No
User:	tpadmin *
Password:	***** *
Host:	tp-c1-ctm-1 *
Registration Status:	Registration completed.

* = Required Fields

Buttons: Apply, Reset

System Status:

Active Meetings:	0
Segments Used:	0
Errors:	0
Warnings:	0
Status:	✓



Note

The CTS Manager and CTMS currently do not utilize longitude and latitude coordinates configured in [Figure 10-16](#).

Quality of Service

The DSCP markings for both the media (audio and video) as well as the signaling from the CTMS are controlled via the QoS Settings screen within the CTMS, shown in [Figure 10-17](#).

Figure 10-17 QoS Settings for the CTMS



Cisco recommends the DSCP marking for audio and video media from the CTMS be CS4, based on RFC 4594. Likewise the recommended DSCP marking for signaling from the CTMS is CS3. This is consistent with recommendations for CTS endpoints defined within CUCM for point-to-point TelePresence meetings.

Meeting Security

Providing secure multipoint meetings is a requirement for any technology—audio, video conferencing, or Cisco TelePresence. Although it is important to provide users with options for multipoint meetings, it is also important that meetings provide the appropriate level of security for participants.

Meeting security can be broken down into three broad categories:

- Administrative access control
- Meeting access control
- Meeting confidentiality

Administrative Access Control

The CTMS provides administrative access control via the Access Management screen as shown in [Figure 10-18](#).

Figure 10-18 Administrative Roles for CTMS Access Control

User-Name	Administrator	Meeting Scheduler	Diagnostic Technician
admin	✓	✓	✓
ctms_scheduler	✗	✓	✗
ctms_technician	✗	✗	✓

Showing 1 - 3 of 3 records

First Previous Next Last Rows per page: 10 New Edit Delete

System Status

Active Meetings: 1
Segments Used: 3
Errors: 0
Warnings: 2
Status: ✓

The CTMS includes three different roles which provide three levels of administrative access control, administrator, meeting scheduler, and diagnostic technician. Multiple instances of each role (i.e., multiple administrators, meeting schedulers, and diagnostic technicians) may be defined within a single CTMS. For each of the three administrative roles, passwords are stored via the local database within the CTMS only. Strong authentication via Radius or TACACS+ to a centralized authentication server is currently not supported.

From a meeting security standpoint, only the administrator and meeting scheduler roles have the ability to schedule meetings. Administrators have full access to the CTMS. Therefore it is recommended to restrict administrator access to a limited set of administrators. Multiple meeting schedulers may be defined within the CTMS. However, it should be noted that resource allocation is not divided among meeting schedulers. In other words, any meeting scheduler has access to all meeting resources within the CTMS. Therefore, meeting security is only as good as the controls which safeguard the access to the CTMS administrator and meeting scheduler roles.

Meeting Access Control

The CTMS provides three types of meetings, as described above, each providing different levels of access control to join the meeting. Scheduled and ad hoc meetings are inherently secure, not allowing uninvited Cisco TelePresence rooms to randomly dial into a meeting. Cisco TelePresence rooms may be added to these meetings once the call is in progress, but only by a meeting scheduler or administrator through the CTMS administrative GUI.

Static meetings are always available, non-scheduled dial-in meetings that are less secure than scheduled and ad hoc meetings, since any room can dial into them at any time. Configuring a hosted static meeting adds a measure of security, in that the host room must at least attend in order for the conference to begin.

The following additional features can be added to multipoint meetings in order to provide an additional level of meeting access control. The Video Announce feature forces each new site entering the conference to be visible to the other sites for approximately two seconds. This minimizes the chance of an unauthorized site “lurking” on the multipoint meeting. Selecting a static meeting in which the number

of rooms defined for the static meeting, via the Maximum Rooms parameter, matches the number of rooms for the conference, minimizes the chance of an unauthorized site being able to access the call. Finally having a meeting scheduler or administrator implement the Lock Meeting feature guarantees that no unauthorized sites can join a meeting once it is ongoing. Note that all meeting features listed above (except Maximum Rooms, which is only available for static and hosted static meetings) may be applied to static, scheduled, or ad hoc meetings.

Meeting Confidentiality

Encryption of the media between CTS endpoints in a point-to-point call has been supported as of CTS version 1.2. Therefore, the recommended method of providing meeting confidentiality for **point-to-point** TelePresence meetings is to enable meeting encryption via the CTS endpoints themselves. However, as of software version 1.1, the CTMS does not support encryption. Therefore, the only method of providing meeting confidentiality for **multipoint** TelePresence meetings is to provide encryption between sites via network components (IPsec encryption). Future versions of CTMS software will include the ability of the CTMS to support encryption.

Meeting Management

Multipoint meeting management is available through the CTMS Web interface, allowing administrators to monitor and or modify active multipoint meetings. The active meetings page allows the administrator to view active meeting details and participants. Administrators can also make the following changes to any active multipoint meeting from the active meeting page:

- Add a room
- Remove a room
- Change switching policy
- Enable or disable video announce
- Enable or disable VIP mode

Figure 10-19 shows the meetings management page of the CTMS.

Figure 10-19 Meeting Management

Meeting ID: 9076374188

Meeting Description: test

Meeting Type: IMMEDIATE

Room1: 1005 Delete Room

Room2: 1003 Delete Room

Add Room(s): * Enter multiple numbers separated by Carriage Return (the ENTER Key).

Switching Policy: Site Segment

Video Announce: Yes No

Quality: Highest Detail, Best Motion: 1080p

VIP Mode: Yes No

Save Close

For centralized multipoint deployments with one CTMS, locating an active multipoint meeting is not an issue. However, for large deployments with multiple CTMS devices, locating a specific multipoint meeting can be a challenge. In any scheduled meeting environment it is recommended that CTS Manager be used to view and manage all meetings. CTS Manager provides a central location to view all past, active, and pending meetings. Figure 10-20 shows the meetings page in CTS Manager.

Figure 10-20 Meetings Page in CTS Manager

Start Time (GMT -7.0)	End Time (GMT -7.0)	Status	Room	Scheduler	Subject
03/12/2008 04:00 PM	03/12/2008 04:30 PM		CTS1 CTS2	aglowack@t...	test
03/13/2008 10:00 AM	03/13/2008 10:30 AM		CTS2 CTS1	aglowack@t...	[Empty Subject]
04/07/2008 03:00 PM	04/07/2008 03:30 PM		CTS1 CTS2	aglowack@t...	Staff Meeting
04/07/2008 04:30 PM	04/07/2008 05:30 PM		CTS2 CTS1	aglowack@t...	Customer Meeting
04/09/2008 08:00 AM	04/09/2008 08:30 AM		CTS1 CTS2	aglowack@t...	Support call
04/09/2008 10:00 AM	04/09/2008 11:30 AM		CTS2 CTS1	aglowack@t...	Support Call

224586

Using the meetings page, you can locate the CTMS a specific multipoint meeting is hosted on by selecting the meeting and clicking Details. In the meeting details you see the CTMS hosting the meeting at the bottom of the page. Proceed to the MCU devices page on CTS Manager, click on the appropriate CTMS, and you are linked directly to the CTMS Web interface. From the active meetings page you can monitor or modify an active meeting.

From the MCU devices page you can also view all past, active, or pending meetings for any CTMS device by selecting the device and clicking View Meetings. From this page you can also export CDR information for all meetings hosted on each CTMS by using the Export Data function.



CHAPTER 11

Cisco Multipoint Technology and Design Details

Audio and Video Flows In A Multipoint TelePresence Design

This section discusses in detail the audio and video flows within a multipoint TelePresence virtual meeting. The information provided within this section can be used by the network design engineer to correctly provision bandwidth to support multipoint TelePresence deployments.

Flow Control Overview

To help control bandwidth use during multipoint meetings, a new flow control feature has been implemented between the CTMS and CTS systems. This feature provides the ability for inactive table segments in a multipoint meeting to stop transmitting video, lowering overall bandwidth utilization.

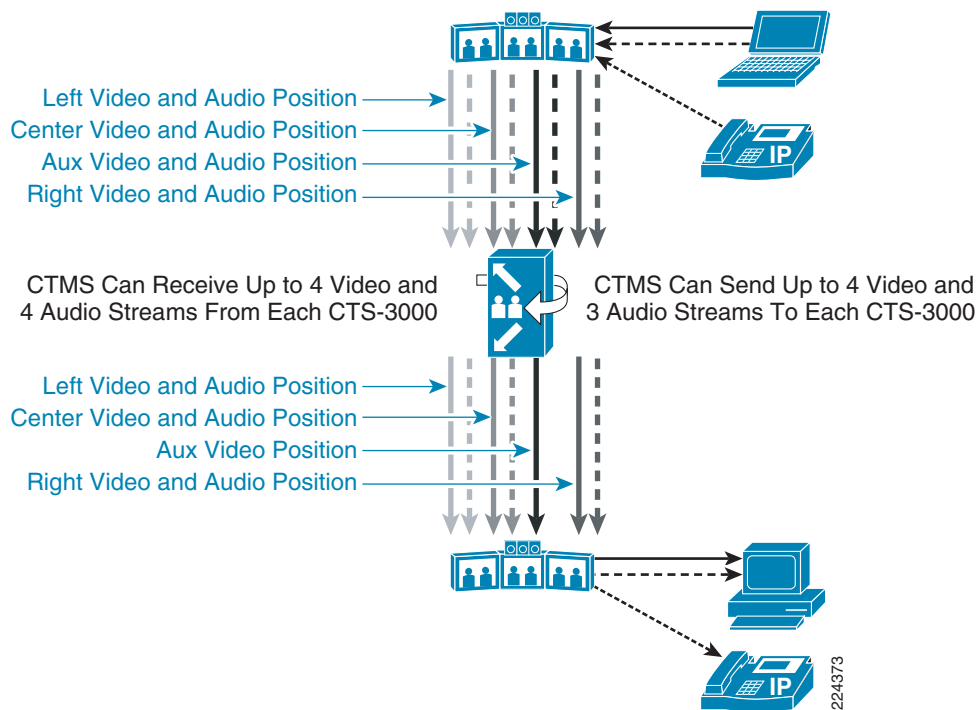
After the multipoint meeting is initiated and the active table segments have been established, CTMS instructs CTS endpoints to stop transmitting video for table segments that are not currently being displayed. Audio continues to be sent from all table segments and used by the CTMS to determine when an inactive table segment becomes active. At that point, the CTMS instructs the CTS system to start transmitting video again for the newly active table segment. This process is continued throughout the meeting, helping reduce overall bandwidth consumption for the multipoint meeting. It should be noted that the CTMS implements a hold-down timer of approximately two seconds to keep the video from flapping due to temporary noise within the room. The flow control feature does not introduce any perceptible delay on top of the hold down time to participants within the meeting.

Audio and Video Positions

CTS endpoints are capable of both sending and receiving multiple audio and video streams. When CTS endpoints join a multipoint call, they first exchange Real Time Control Protocol (RTCP) packets. Successful exchange of these packets indicates the opposite endpoint is a Cisco TelePresence device, capable of supporting various Cisco extensions. Among other things, these extensions are used to determine the number of audio and video channels each TelePresence endpoint may send and receive.

Audio and video streams are sent and received based on their position within the CTS endpoint. [Figure 11-1](#) shows this for a multipoint call consisting of CTS-3000s.

Figure 11-1 Audio and Video Stream Positions with CTS-3000s



Each CTS-3000 can transmit up to four audio streams and four video streams from the left, center, right, and auxiliary positions. These correspond to the left, center, and right cameras and microphones, as well as the auxiliary input. The auxiliary video input can be used by a PC for slide show presentations. The auxiliary audio input is shared between audio from the PC that accompanies the slide show presentation and audio from an audio-only participant, such as an IP phone add-on.

**Note**

With the CTS-3000 all microphones physically connect to the center codec, even though the audio positions are referred to as center, left, or right.

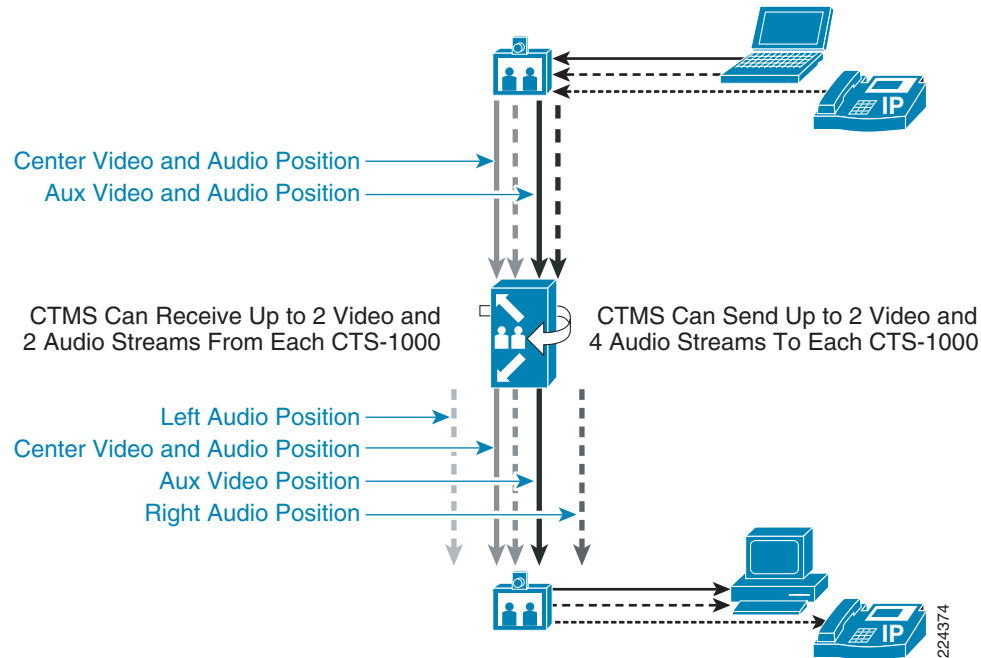
The CTMS can transmit up to four video streams, corresponding to the left, center, and right plasma displays of the CTS-3000, as well as either a projector or monitor for slide show presentations connected to the auxiliary video output. The CTMS can only transmit up to three audio streams, corresponding to the left, center, and right speaker positions of the CTS-3000. Audio sent by an originating CTS-3000 toward the auxiliary position is redirected to one of the three speaker positions of the destination CTS-3000 by the CTMS. The CTMS chooses the three loudest audio streams to send to the remote CTS-3000 when there are more than three streams with audio energy.

**Note**

The number of audio streams sent in a multipoint call is different than in a point-to-point call in which four audio streams can be sent and received by each CTS-3000.

Figure 11-2 shows the audio and video positions for a multipoint call consisting of CTS-1000s.

Figure 11-2 Audio and Video Stream Positions with CTS-1000s



Each CTS-1000 can transmit up to two audio streams and two video streams from the center and auxiliary positions. These correspond to the single camera and microphone of the CTS-1000, as well as the auxiliary input. However, the CTMS can still transmit up to three audio streams, corresponding to the left, center, and right speaker positions of the CTS-1000, even though the CTS-1000 only has a single center speaker. The CTS-1000 mixes the audio from each of the three positions to play out on its single speaker. The CTS-1000 can only receive up to two video streams, corresponding to the center plasma display and either a projector or monitor for slide show presentations connected to the auxiliary video output.

Audio to the CTMS in a Multipoint TelePresence Meeting

Audio from all TelePresence conference participants is always sent to the CTMS, regardless of whether the site has any audio energy (someone is speaking) or not. In other words, silence suppression of audio packets is not implemented within CTS units. Each microphone of a CTS unit transmits a single 64 Kbps RTP/AVC (IETF RFC 3551) audio stream using a separate RTP SSRC. A second 64 Kbps audio stream can be sent from either the auxiliary audio input or via the audio add-on feature, which can be used to add either a single phone or an audio conferencing bridge into the TelePresence virtual meeting.

Audio packet sizes average approximately 220 bytes in size including network headers and are sent every 20 ms. The payload size is approximately 160 bytes. Therefore the network overhead used for these calculations is approximately 27.27%. Each inbound audio stream generates approximately 88 Kbps inbound to the CTMS across an Ethernet segment.



Note

This has been confirmed through data traces taken by ESE. It includes a 20 Byte IP header, 8 Byte UDP header, 12 Byte RTP header, and 14 byte Ethernet header. When calculating the amount of bandwidth across the WAN, the Ethernet header overhead must be replaced with the appropriate Layer 2 WAN header.

Calculating the Amount of Audio Traffic to the CTMS

The number of audio streams inbound to the CTMS in a single multipoint meeting can be calculated by the following equations:

$(N + (3 * M))$ when the auxiliary audio input and audio-only add-on is not used

Or:

$(N + (3 * M) + P)$ when the auxiliary audio input and/or an audio-only phone is added on

Where N = the number of CTS-1000 endpoints in the call, M = the number of CTS-3000 endpoints in the TelePresence call, and P is the number of auxiliary audio inputs and audio-only phones added-on to the meeting.

The total amount of audio traffic inbound to the CTMS in a single multipoint meeting can be calculated by simply multiplying the audio bandwidth per call to the equations above to yield the following:

$88 \text{ Kbps} * (N + (3 * M))$ when the auxiliary audio input audio-only add-on is not used

Or:

$88 \text{ Kbps} * (N + (3 * M) + P)$ when the auxiliary audio input and/or an audio-only phone is added on

Note that only one device in a multipoint call can function as “presenter” for auxiliary video (i.e., PowerPoint slides) and audio input. Also, multiple audio-only phones can be bridged on through separate sites. However, if multiple audio-only devices need to be added into a TelePresence multipoint meeting, it is more effective to add an audio bridge, rather than have multiple sites add individual audio-only phones.

As an example, a 6-site multipoint CTS-1000 call in which the auxiliary audio input is not used and no audio-only phones are bridged onto the TelePresence meeting has an estimated inbound audio rate to the CTMS of the following:

$88 \text{ Kbps per audio stream} * 6 \text{ CTS-1000s} = 528 \text{ Kbps toward the CTMS}$

The total amount of audio traffic inbound to the CTMS from multiple multipoint meetings can be calculated by simply summing the traffic from individual meetings. Extending the example above, if the CTMS is currently supporting one 6-site multipoint call with CTS-1000s only, and one 3-site multipoint call with CTS-3000s and an audio conference bridge added on, the total amount of inbound audio could be calculated as:

$(88 \text{ Kbps} * 6 \text{ CTS-1000s}) + (88 \text{ Kbps} * (3 * 3 \text{ CTS-3000s} + 1 \text{ Audio Conf. Bridge}))$

$528 \text{ Kbps} + 880 \text{ Kbps} = 1.408 \text{ Mbps}$

The maximum number of inbound audio streams to a CTMS can be estimated based upon the maximum number of table segments supported by the CTMS. Assuming 48 CTS-1000s in 16 separate three-party multipoint calls, with each site having an audio-only phone bridged onto the TelePresence meeting, the total amount of inbound audio traffic to the CTMS can be estimated as:

$(3 \text{ CTS-1000s} + 3 \text{ audio-only add-on phones}) = 6 \text{ inbound audio streams per multipoint call}$

$6 \text{ inbound audio streams} * 16 \text{ multipoint calls} = 96 \text{ inbound audio streams to the CTMS}$

The maximum amount of audio traffic inbound to the CTMS can be calculated by simply multiplying the audio bandwidth per call with the equations above to yield the following:

$88 \text{ Kbps} * 96 \text{ inbound audio streams} = 8.448 \text{ Mbps}$

Therefore, the network would need to be able to support approximately 8.5 Mbps of inbound audio to the CTMS. Since the audio traffic is marked with the same DSCP marking (recommended by Cisco to be CS4) as the video traffic in a TelePresence meeting, this amount of audio may be relatively small compared to the amount of inbound video traffic to the CTMS. Video traffic calculations are discussed in other sections.

Audio From the CTMS in a Multipoint Meeting

Audio sent from the CTMS is somewhat more complex than audio sent to the CTMS. Since audio is continuously sent from each CTS unit to the CTMS, there can be multiple simultaneous speakers in a TelePresence multipoint conference. The CTMS determines which video stream is replicated and sent to the other endpoints, based on which speakers are talking at the moment and who is the loudest. This is signaled by way of the voice activity confidence metric within every voice packet sent from every CTS endpoint. Within the initial RTCP packet exchange which occurs immediately after the CTS endpoint establishes a connection with the CTMS, the CTS endpoint advertises its capability to send a voice activity confidence metric within voice packets. The voice activity confidence metric is an estimation of the amount of audio energy contained within the voice packet.



Note

The CTMS does not advertise the ability to send a voice activity confidence metric to CTS endpoints, nor does it include a voice activity confidence metric within voice packets sent to CTS endpoints.

Audio is replicated when sufficient audio energy is detected within the voice packets from each CTS endpoint, a feature known as Voice Activity Switching. The CTMS replicates up to three audio streams, each destined to a particular audio position of the CTS endpoint, left, center, or right. Each CTS endpoint mixes the inbound audio to be sent to its audio speaker(s). The audio data rate outbound from the CTMS onto the network is therefore somewhat variable, based on the number of simultaneous speakers in the virtual meeting.

Calculating the Amount of Audio Traffic from the CTMS

The number of audio streams outbound from the CTMS in a single multipoint meeting varies based upon how many speakers are talking. However, an estimate of the maximum number of audio streams outbound from the CTMS for a given multipoint meeting can be calculated with the following equation:

$$3 \text{ Audio Streams} * \text{Number of CTS Endpoints in the multipoint call}$$

The total amount of audio traffic outbound from the CTMS in a single multipoint meeting can be calculated by simply multiplying the audio bandwidth per call with the equation above to yield the following:

$$88 \text{ Kbps} * 3 \text{ Audio Streams} * \text{Number of CTS Endpoints in the multipoint call}$$

For example, in a 6-site multipoint CTS-1000 call in which all sites are receiving the maximum of three audio streams, the CTMS would be sending approximately the following:

$$3 \text{ Audio Streams} * 6 \text{ CTS-1000s} = 18 \text{ Audio Streams}$$

$$88 \text{ Kbps per Audio Stream} * 18 \text{ Audio Streams} = 1.58 \text{ Mbps of audio traffic}$$

The total amount of audio traffic outbound from the CTMS from multiple multipoint meetings can again be calculated by simply summing the traffic from individual meetings.

Finally, the maximum number of outbound audio from a CTMS can also be estimated based upon the maximum number of audio segments supported by the CTMS. Assuming 48 CTS-1000s in 16 separate three-party multipoint calls, each receiving the maximum of three audio streams, the total number of outbound audio streams from the CTMS can be estimated as:

$$3 \text{ Audio Streams} * 3 \text{ CTS-1000s per Call} * 16 \text{ Separate Calls} = 144 \text{ Audio Streams}$$

The maximum amount of audio traffic outbound from the CTMS can be estimated by simply multiplying the audio bandwidth per call with the equation above to yield the following:

$$88 \text{ Kbps} * 144 \text{ Audio Streams} = 12.672 \text{ Mbps of audio outbound from the CTMS}$$

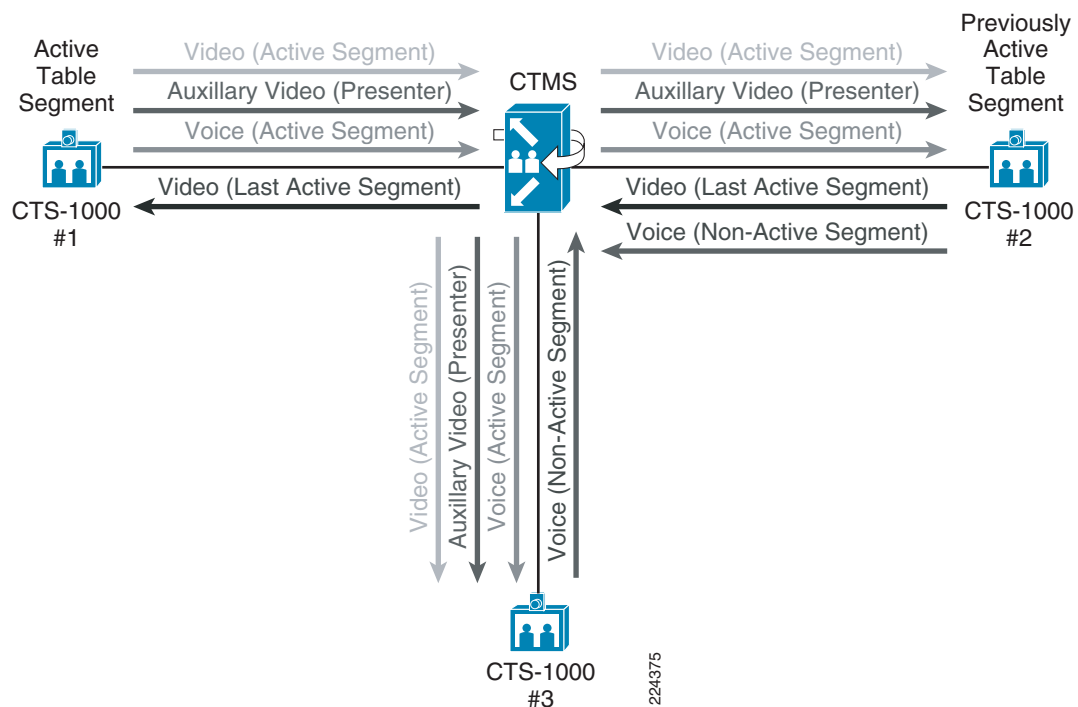
Therefore, the network would need to be able to support approximately 12.7 Mbps of outbound audio from the CTMS.

Comparing the amount of audio sent to the CTMS with the amount of audio sent from the CTMS indicates that considerably more audio is typically sent outbound from the CTMS than is received by the CTMS during a TelePresence meeting. This traffic pattern is indicative of a multipoint meeting in which multiple audio streams have to be replicated and sent to CTS endpoints which can each simultaneously receive multiple audio streams. This is one of the major network differences between multiple point-to-point TelePresence meetings and multipoint TelePresence meetings.

Video in a Multipoint TelePresence Meeting

Unlike audio, video is not continuously transmitted from each CTS endpoint to the CTMS. Instead, the CTMS signals which endpoint should send its video. The CTMS determines which video stream to present to TelePresence meeting participants based on which speaker is currently talking or which speaker is talking the loudest if multiple speakers are talking simultaneously, also known as the active site or active segment. Figure 11-3 shows an example of this in a three-site CTS-1000 TelePresence call.

Figure 11-3 Video Flows in a 3-Site TelePresence Call



In order to switch the video, the CTMS determines which site is the active segment based upon the value of the voice activity confidence metric transmitted within voice packets from each site. Note that for multipoint calls which include CTS-3000s using speaker switching, there can be multiple active segments.

In the example above, CTS-1000 #1 is the active segment and slide presenter. Video from CTS-1000 #1 is therefore replicated on a packet-by-packet basis by the CTMS and sent to CTS-1000 #2 and CTS-1000 #3. However, the display of CTS-1000 #1 needs to continue showing video from the previous active

segment. In the example above, the last active segment was CTS-1000 #2. Therefore, video from CTS-1000 #2 continues to be sent to the CTMS, where it is replicated on a packet-by-packet basis and sent to CTMS-1000 #1.

Camera Video Input

Video from the camera inputs is transmitted via H.264 at 30 frames/second using a separate RTP SSRC for each camera. Unlike audio, TelePresence video can be sent at different overall bit rates based upon the quality configuration of the CTS-endpoints and the CTMS. The following discussion is based upon the multipoint call configured for 1080p Best quality. Video rates can burst up to $4 \text{ Mbps} * 110\% = 4.4 \text{ Mbps}$ per camera with this video setting. Since average video packet sizes are 1,100 bytes, network overhead can add an additional 4.91% overhead for a total of approximately 4.616 Mbps per video stream.

Auxiliary Video Input

Cisco TelePresence currently supports two frame rates for auxiliary video input, low speed auxiliary video input at 5 frames per second and high speed auxiliary video input at 30 frames per second. High speed auxiliary video input requires a separate codec be added to existing CTS endpoints and is not covered in this document.

Low speed auxiliary video is transmitted via H.264 at 5 frames/second using a separate RTP SSRC again. The maximum data rate can burst up to approximately $500 \text{ Kbps} * 110\% = 550 \text{ Kbps}$. Since average video packet sizes are approximately 1,100 bytes, network overhead can add another 4.91% overhead for a total of 577 Kbps for the low speed auxiliary video stream.



Note

These video overhead calculations include a 20 Byte IP header, 8 Byte UDP header, 12 Byte RTP header, and 14 byte Ethernet header. When calculating the amount of bandwidth across the WAN, the Ethernet header overhead must be replaced with the appropriate Layer 2 WAN header.

Calculating the Amount of Video Traffic to the CTMS

The total number of inbound video streams to the CTMS varies based upon the type of CTS units involved in the TelePresence meeting.

Meetings with CTS-1000s Only

For meetings which involve only CTS-1000s, the following equations hold for meetings with and without the use of the auxiliary video input, regardless of the number of CTS-1000 units involved in the multipoint call.

2 Inbound Camera Video Streams without Auxiliary Video Input Stream

Or:

2 Inbound Camera Video Streams + 1 Inbound Auxiliary Video Input Stream

The camera video input streams correspond to the active site and the last active site. There will only be a single additional inbound auxiliary video stream if any of the CTS-1000 units is functioning as a presenter (i.e., using the auxiliary video input for a PowerPoint presentation).

**Note**

When multiple devices connect to the auxiliary video and audio input, the last device connected is the presenter.

Therefore the following equations can be used to estimate the total amount of inbound video traffic to the CTMS by multiplying the video stream rates by the number of streams:

$$4.616 \text{ Mbps} * 2 \text{ Inbound Video Streams} = 9.232 \text{ Mbps without auxiliary video input}$$

$$4.616 \text{ Mbps} * 2 \text{ inbound video streams} + 577 \text{ Kbps} = 9.809 \text{ Mbps with auxiliary video input}$$

**Note**

In order to determine the total amount of video with high-speed auxiliary video input, simply substitute 577 Kbps with 4.616 Mbps in all the equations discussed in this section.

The total number of inbound video streams to the CTMS from multiple TelePresence meetings involving only CTS-1000s can be found by multiplying the equations above by the number of simultaneous TelePresence meetings supported. For example, the total number of inbound video streams to a CTMS which is currently supporting two 8-site meetings and four 3-site meetings, both with presenters showing PowerPoint slides, can be calculated as:

$$6 \text{ total meetings} * (2 \text{ Inbound Camera Video Streams} + 1 \text{ Inbound Auxiliary Video Stream}) = 12 \text{ Inbound Camera Video Streams} + 6 \text{ Inbound Auxiliary Video Streams}$$

The total amount of inbound video traffic to the CTMS can be estimated by simply multiplying the video rates by the number of meetings:

$$6 \text{ meetings} * 9.809 \text{ Mbps Per Meeting with Auxiliary Video Input} = 58.854 \text{ Mbps}$$

The maximum amount of inbound video to a CTMS supporting only CTS-1000s can be estimated based upon the maximum number of video segments supported by the CTMS. Assuming 48 CTS-1000s in 16 separate three-party multipoint calls, each using the auxiliary video input, the total amount of inbound video traffic to the CTMS can be estimated as:

$$16 \text{ meetings} * 9.809 \text{ Mbps Per Meeting with Auxiliary Video Input} = 156.944 \text{ Mbps}$$

Therefore, the network would need to be able to support approximately 157 Mbps of inbound video to the CTMS.

Meetings with CTS-3000s and CTS-3200s Only

For meetings which involve only CTS-3000s and CTS-3200s, there are six inbound camera video streams, regardless of the number of CTS-3000 and CTS-3200 units involved in the call and regardless of whether speaker-switching or room-switching is implemented within the meetings. These video streams correspond to the active segments and the last active segments. There is also a single additional inbound auxiliary video stream if any of the CTS-3000 or CTS-3200 units is functioning as a presenter (i.e., using the auxiliary video input for a PowerPoint presentation).

Therefore the following equations can be used to estimate the total amount of inbound video traffic to the CTMS, again assuming low-speed auxiliary video:

$$4.616 \text{ Mbps} * 6 \text{ Inbound Video Streams} = 27.696 \text{ Mbps without Auxiliary Video Input}$$

$$4.616 \text{ Mbps} * 6 \text{ inbound Video Streams} + 577 \text{ Kbps} = 28.273 \text{ Mbps with Auxiliary Video Input}$$

The total amount of inbound video to the CTMS from multiple TelePresence meetings involving only CTS-3000s and CTS-3200s can be found by multiplying the equations above by the number of simultaneous TelePresence meetings supported. For example, the total amount of inbound video traffic to a CTMS which is currently supporting two 4-site meetings and two 3-site meetings, both with presenters showing PowerPoint slides, can be calculated as:

$$4 \text{ Total Meetings} * 28.273 \text{ Mbps Per Meeting with Auxiliary Video Input} = 113.092 \text{ Mbps}$$

The maximum amount of inbound video to a CTMS supporting only CTS-3000s and CTS-3200s can be estimated based upon the maximum number of video segments supported by the CTMS. Assuming 16 CTS-3000s in four separate 3-site multipoint calls and one 4-site multipoint call, each using the auxiliary video input, the total amount of inbound video traffic to the CTMS can be estimated as:

$$5 \text{ Total Meetings} * 28.273 \text{ Mbps Per Meeting with Auxiliary Video Input} = 141.365 \text{ Mbps}$$

Therefore, the network would need to be able to support approximately 141 Mbps of inbound video traffic to the CTMS.

Meetings with Combinations of CTS-1000s, CTS-3000s, and CTS-3200s

For meetings which involve one CTS-3000 or CTS-3200 and two or more CTS-1000s, or meetings which involve two or more CTS-3000s or CTS-3200s, and any number of CTS-1000s, there are always six inbound camera video streams. These video streams correspond to the active segments and the last active segments. There is also a single additional inbound auxiliary video stream if any of the CTS units is functioning as a presenter (i.e., using the auxiliary video input for a PowerPoint presentation). Therefore the same equations from the previous section regarding CTS-3000s and CTS-3200s only apply to mixed meetings of CTS-1000s, CTS-3000s, and CTS-3200s.

The one exception is when there is one CTS-3000 or CTS-3200 and only two CTS-1000s. In this case, there are a total of five inbound camera video streams. It should be noted that there are not enough video streams to fill all the CTS-3000 displays in such a meeting. In other words one CTS-3000 screen is blank.



Note

It is assumed that meetings with one CTS-3000 or CTS-3200 and one CTS-1000 do not require a CTMS, and are not considered a multipoint meeting, although it is possible to hold a two-site multipoint meeting.

Calculating the Amount of Video Traffic From the CTMS

The total number of outbound video streams from the CTMS to the CTS units equals the number of video table segments in the multipoint call, where each CTS-1000 counts as one table segment and each CTS-3000 or CTS-3200 counts as three table segments. In addition, if the auxiliary video input is being used in the TelePresence meeting, an additional amount of video up to 577 Kbps times the number of endpoints is transmitted by the CTMS. Again this assumes low-speed auxiliary video only.

The following equation can be used to estimate the amount of outbound video from the CTMS:

$$((N + (3 * M)) * 4.616 \text{ Mbps Per Video Stream without Auxiliary Video Input}$$

$$((N + (3 * M)) * 4.616 \text{ Mbps} + (N + M - 1) * 577 \text{ Kbps with Auxiliary Video Input}$$

Where N is the number of CTS-1000s in the call and M is the number of CTS-3000s in the call.

For example, in a single multipoint call with 3 CTS-3000s and 5 CTS-1000s, along with an auxiliary video stream from a presentation, the estimated amount of outbound video traffic from the CTMS would be:

$$((5 + (3 * 3)) * 4.616 \text{ Mbps} + (5 + 3 - 1) * 577 \text{ Kbps} = 68.663 \text{ Mbps}$$

The total amount of video traffic outbound from the CTMS from multiple multipoint meetings can be calculated by simply summing the traffic from individual meetings.

Extending the example above, if the CTMS is currently supporting one multipoint call with 3 CTS-3000s and 5 CTS-1000s, along with a second call consisting of four CTS-3000s, both using the auxiliary video input for PowerPoint presentations, the total amount of video traffic outbound from the CTMS can be estimated as:

$$[((5 + (3 * 3)) * 4.616 \text{ Mbps} + (5 + 3 - 1) * 577 \text{ Kbps}) + [(4 * 3) * 4.616 \text{ Mbps} + (4 - 1) * 577 \text{ Kbps}]] = 125.786 \text{ Mbps}$$

The maximum amount of outbound video from a CTMS can be estimated based upon the maximum number of video segments supported by the CTMS. Assuming 48 CTS-1000s in one large multipoint call, using the auxiliary video input for a PowerPoint presentation, the total amount of outbound video traffic from the CTMS can be estimated as:

$$[48 * 4.616 \text{ Mbps} + (48 - 1) * 577 \text{ Kbps}] = 248.687 \text{ Mbps of Video Outbound From the CTMS}$$

Therefore, the network would need to be able to support approximately 249 Mbps of outbound video from the CTMS.

Total Traffic to and from the CTMS

The total amount of traffic to and from the CTMS for a given meeting or set of meetings can be found by summing the amount of audio and video presented in the previous sections. For example, the maximum amount of traffic from the CTMS can be estimated as:

$$249 \text{ Mbps of Video} + 13 \text{ Mbps of Audio} = 262 \text{ Mbps}$$



Note

One could simply estimate the total traffic from the CTMS by using the estimated 5.5 Mbps per CTS-1000 and 15 Mbps per CTS-3000 or CTS-3200 and multiplying by the number of devices supported. For example 48 CTS-1000s would yield a total amount of traffic of $48 * 5.5 \text{ Mbps} = 264 \text{ Mbps}$. Likewise, 16 CTS-3000s or CTS-3200s would yield a total amount of traffic of $16 * 15 \text{ Mbps} = 240 \text{ Mbps}$.

Simply estimating 15 Mbps per CTS-3000 or CTS-3200 and 5.5 Mbps per CTS-1000 or CTS-500 provides reasonably accurate numbers for traffic outbound from the CTMS across the network. However, due to the asymmetric nature of multipoint TelePresence, they do not accurately reflect the amount of traffic inbound to the CTMS across the network. Further, with the addition of high speed auxiliary video input, a single traffic rate utilized for CTS endpoints in calculating bandwidth utilization will become increasingly inaccurate depending upon whether the auxiliary video input is in use or not in use. For this reason, the detailed explanation within this document has been provided.

Video Switchover Delay

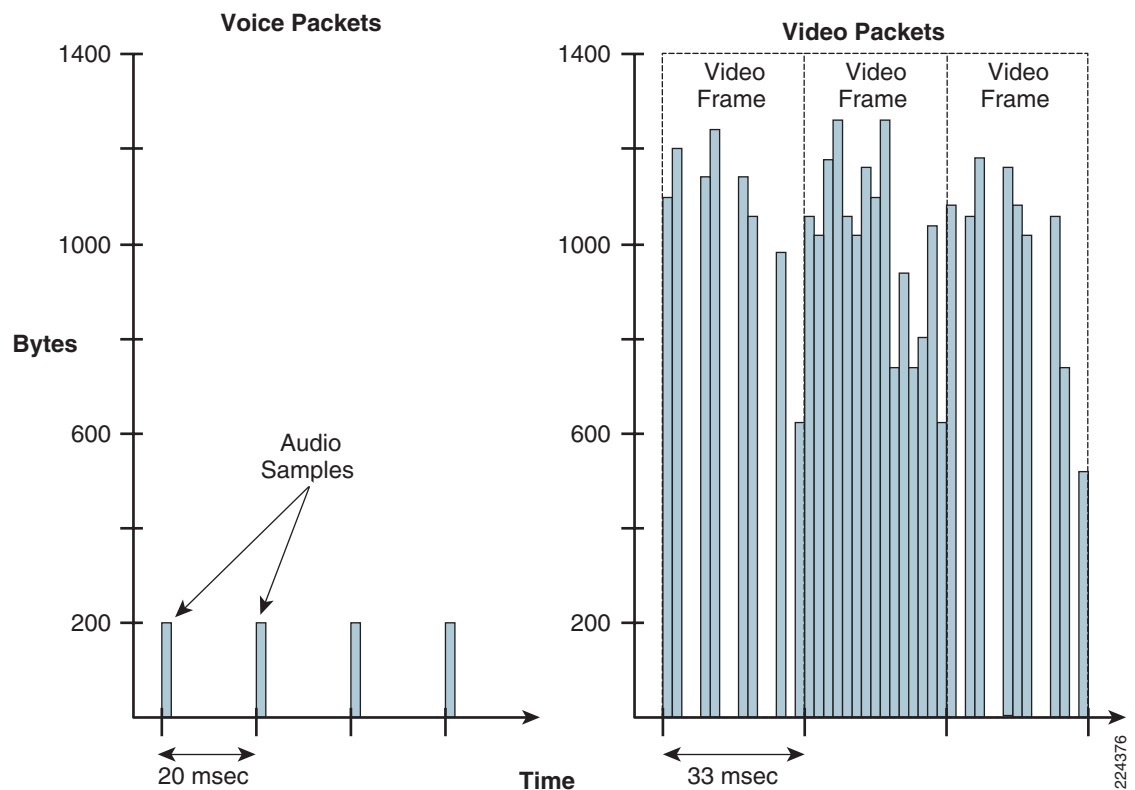
The CTMS does not immediately signal CTS endpoints to send video upon seeing packets with audio energy. This could cause unnecessary flapping of video within the conference call and generate additional burstiness on the network. Instead the CTMS implements a hold down timer before signaling the new active segment to begin sending video. The hold down timer is designed to ensure that the new active segment is indeed speaking, and not just a random noise. There is also a short period of time between when the new active segment has been notified to start sending video and when it begins transmitting video, which is then switched by the CTMS. Any participants within a TelePresence multipoint conference should be aware that they may need to talk for approximately two seconds before

the video switches over. They should be particularly aware of this when taking a roll-call of participants at the beginning of the TelePresence meeting, in order for their faces to be seen on the video. It should be noted that the audio is never interrupted or delayed.

Overview of TelePresence Video on the Network

Since many network configuration guidelines around the deployment of TelePresence, both within the campus and to the branch, are based on the specific behavior of video on the network, a brief overview of TelePresence video as it appears on the network is presented in this section. Figure 11-4 shows a sample comparison of voice and video traffic as it appears on a network.

Figure 11-4 Comparison of Voice and Video on the Network

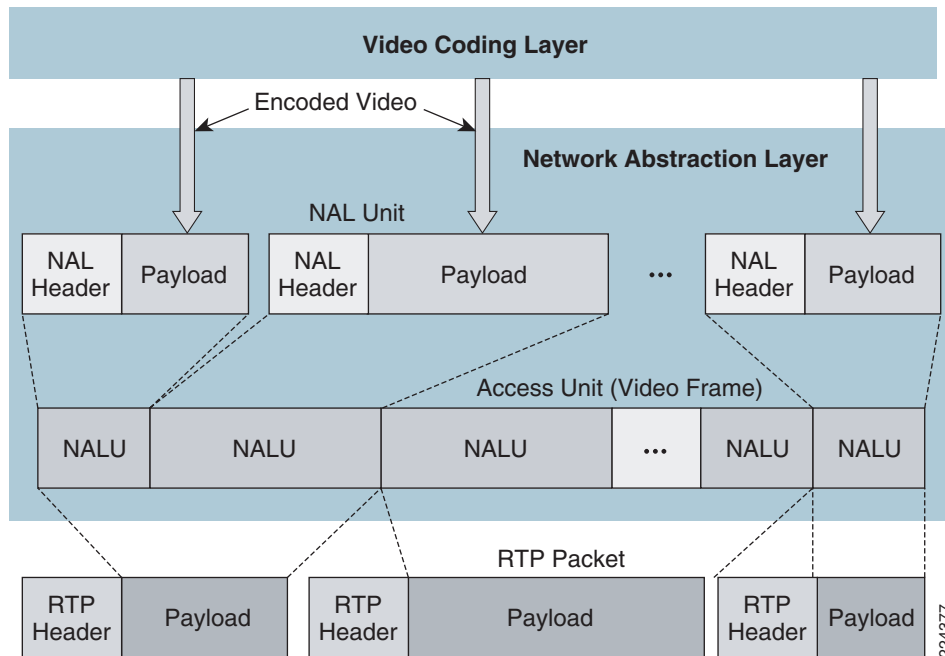


Voice on the network appears as a series of packets, spaced at regular intervals (in the case of Cisco TelePresence, every 20 ms), each containing an encoded sample of the audio. Each voice packet is basically independent of the other packets. In other words, if one voice packet is discarded or lost in the network, it does not affect the next voice packet. The sizes of the voice packets are fairly consistent, averaging slightly over 200 bytes in size. Therefore, the overall characteristic of voice is a constant bit rate stream.

Unlike voice, the overall characteristic of video in general, including TelePresence, is a somewhat bursty, variable bit rate stream. Video traffic on the network appears as a series of video frames spaced at regular intervals (in the case of TelePresence video, approximately every 33 ms). A frame of video is also referred to as an Access Unit in H.264 terminology. The H.264 standard defines two layers, a Video Coding Layer (VCL) and a Network Abstraction Layer (NAL). The VCL is responsible for encoding the video and the output of the VCL is a string of bits representing the encoded video. The function of the

NAL is to map the string of bits into units which can then be transported across a network infrastructure. Each video frame consists of multiple packets spaced out over the frame interval. Each RTP packet contains one or more NAL Units (NALUs). Each NALU consists of an integer number of bytes of coded video, as shown in Figure 11-5.

Figure 11-5 Mapping TelePresence Video into RTP Packets



RTP Packets within a single video frame, and across multiple frames, are not necessarily independent of each other. In other words, if one packet within a video frame is discarded, it affects the quality of the entire video frame and may possibly affect the quality of other video frames. The sizes of the individual RTP packets within frames vary, depending upon the number of NALUs they carry and the size of the NALUs. Overall, packet sizes average around 1,100 bytes in size. The number of packets per frame also varies considerably based upon how much information is contained within the video frame. This is partially determined by how the video is encoded.

There are basically two types of encoding:

- Intra-frame encoding—Compresses the frame by reducing spatial redundancy within the frame.
- Inter-frame encoding—Uses motion compensation to reduce temporal redundancy across one or more frames.

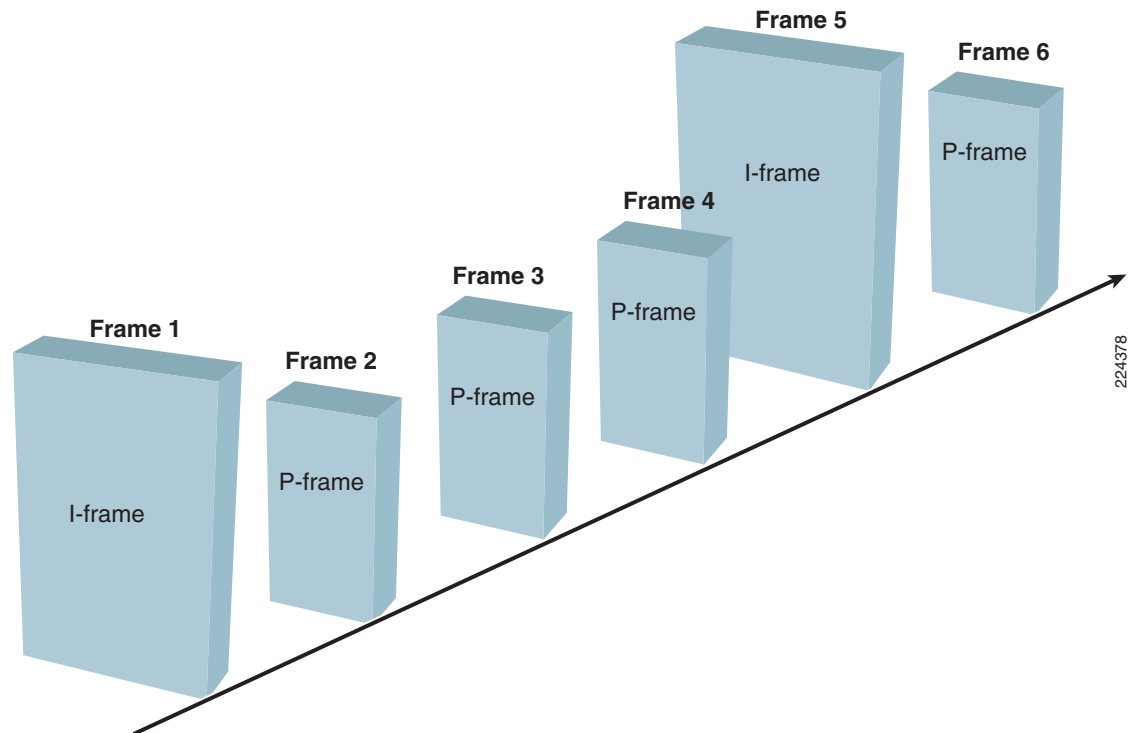
These two types of encoding lead two types of video frames, intra-coded frames (I-frames) and predictive-coded frames (P-frames). A third type of frame, bi-directional predictive-coded frames (B-frames), is currently not utilized by TelePresence.



Note

Coding is actually done at the macroblock layer. An integer number of macroblocks then form a slice, and multiple slices form a frame. Therefore, technically slices are intra-predicted (I-slices) or inter-predicted (P-slices). For simplicity of explanation within this section, these have been abstracted to I-frames and P-frames. A thorough discussion of the H.264 Video Coding Layer is outside the scope of this section.

Figure 11-6 shows an example of how these frames can appear in a video stream.

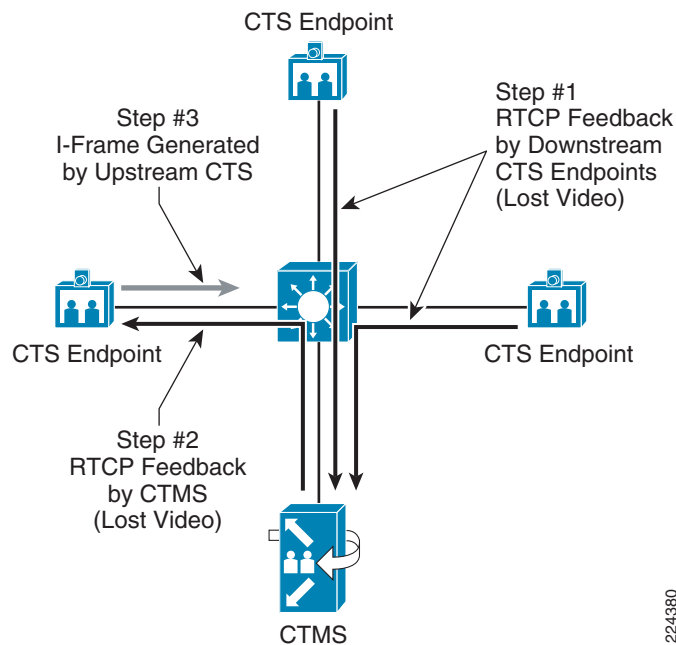
Figure 11-6 Types of Video Frames

I-frames serve as reference points in the video stream. They can also be referred to as an Instantaneous Decoding Refresh (IDR) in H.264 terminology. I-frames can be decoded and displayed without referencing any other frame. Frame 1 in [Figure 11-6](#) is an I-frame. On the other hand, P-frames reference either another P-frame or an I-frame. They require reception of the previous reference frame in order to be decoded correctly. For example Frame 2 in [Figure 11-6](#) references Frame 1. Frame 3 may reference Frame 1 or Frame 2, since a P-frame may reference another P-frame.

Compression of I-frames is typically only moderate, since only spatial redundancy within the frame is eliminated. Therefore, I-frames tend to be much larger in size than P-frames. I-frame sizes up to 64 Kbytes (and approximately 60 individual packets) have been observed with TelePresence endpoints. P-frames have much higher compression since only the difference between the frame and the reference frame is sent. This information is typically sent in the form of motion vectors indicating the relative motion of objects from the reference frame. The size of TelePresence P-frames is dependent upon the amount of motion within the conference call. Under normal motion, TelePresence P-frames tend to average around 13 Kbytes in size and typically consist of about 12 individual packets. Under high motion they can be around 19 Kbytes in size and consist of about 17 individual packets.

From a bandwidth utilization standpoint, much better performance can be achieved by sending I-frames very infrequently. This reduces the burstiness of the video as well as the overall bit rate. However, the side effect is that if part of one video frame is lost (in other words, if a packet is dropped in the network), then multiple video frames which reference it or reference each other may be affected.

In point-to-point TelePresence meetings, I-frames are sent approximately every five minutes. However, waiting for up to five minutes for video to correct itself in the event of a lost packet is not acceptable. Therefore, CTS endpoints include an RTCP-based feedback mechanism by which video receivers continuously update the sender regarding the status of video packets received. When the sender learns that the receiver has lost some video packets, the sender generates an IDR in order to establish a new reference point.

Figure 11-8 Lost Video Feedback in a Multipoint Call

If video has been lost by any of the downstream CTS endpoints, it is reported back to the upstream CTS endpoint by the CTMS after a brief hold-down time. The hold-down time prevents excessive I-frame generation by allowing the CTMS to aggregate reports from all downstream CTS endpoints before informing the upstream CTS endpoint that it needs to send a new I-frame.

The purpose of the feedback mechanism is to ensure that every video endpoint remains synchronized with the source and maintains high video quality. The side effect is that any site which experiences network degradation causes I-frames to be sent and replicated to every site within the multipoint call. In a worst case scenario, the site experiencing network degradation continues to report lost video, resulting in an I-frame “storm” across the network. Therefore, care must be taken to ensure that all sites are provisioned correctly in order to prevent excessive I-frame video traffic.

Deployment Models

There are a number of factors that determine how multipoint is deployed in a production environment. The two biggest factors are the number of CTS endpoints and the geographic location of the CTS endpoints. The number of CTS endpoints determines the number of CTMS devices in the network and the location of the CTS endpoints determines whether the CTMS devices are centralized or distributed.

Centralized Deployment

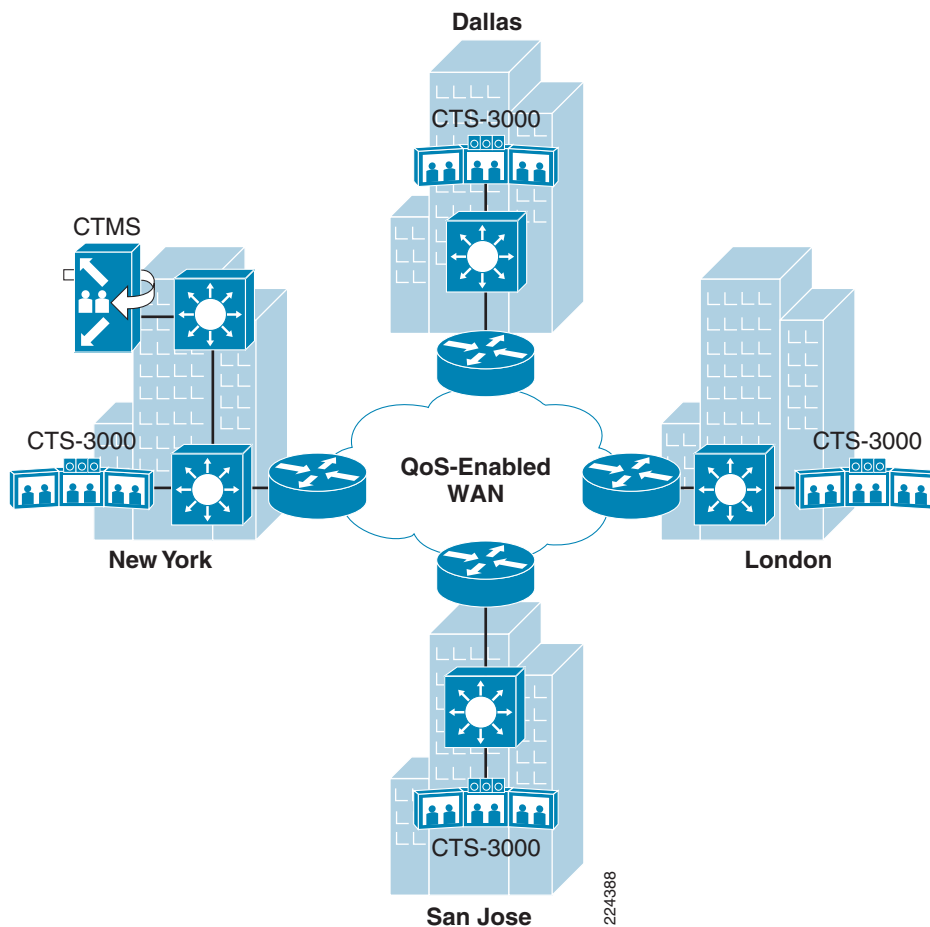
Centralized designs are recommended for Cisco TelePresence deployments with six or fewer CTS units or for larger deployments which cover a limited geographic area.

For centralized deployments, it is recommended that the CTMS be located at a regional or headquarters campus site with the necessary WAN bandwidth available to each of the remote sites, as well as the necessary LAN bandwidth within the campus. It is recommended that the CTMS be centrally located,

based on the geographic location of the CTS rooms, although this may not be entirely possible due to the existing network layout. This prevents unnecessary latency caused by backhauling calls to a site at the far edge of the network.

Figure 11-9 illustrates a small TelePresence deployment with three regional/headquarter campus sites in North America and one site in Europe. In this example the CTMS is placed centrally, located in New York, to minimize latency for multipoint meetings.

Figure 11-9 Centralized Multipoint Design



Deployment Considerations

There are a number of decisions that need to be made prior to deploying a centralized TelePresence multipoint solution.

1. Selecting a site for the multipoint switch—As mentioned above, the multipoint switch should be located at a site that provides end-to-end network latency less than 200ms and provides adequate bandwidth for the number of CTS endpoints on the network. Calculate the required bandwidth for the multipoint site using the calculations from the previous section outlining bandwidth requirements.
2. Supported meeting types—Supporting either scheduled, non-scheduled, or both meeting types is not a concern in most centralized multipoint deployments.

- a. A scheduled only meeting environment requires CTS Manager and provides one-button-to-push meeting access for end users. Meetings can be scheduled by end users or a centralized scheduling group using Microsoft Exchange or IBM Domino. When configuring CTMS resources, maximum segments should be configured for the total available segments plus additional ad hoc segments. It is important to remember that ad hoc resources need to be available even in a scheduled only environment. If no ad hoc resources are configured on the CTMS, there is no way to add non-scheduled CTS endpoints to scheduled meetings (as described in [Multipoint Resources](#) in [Chapter 10](#), “[Cisco TelePresence Multipoint Solution Essentials](#)”). Scheduled segments are calculated based on maximum and ad hoc segment entries. [Figure 11-10](#) illustrates an example of resource allocation for a scheduled only meeting deployment with five CTS-3000s (15 available segments).

Figure 11-10 Resource Allocation Example

IP Settings	Access Settings	QoS Settings	Resource Management	SNMP Settings	Restart CTMS
Maximum Segments:	21	*			
Adhoc Segments:	6	*			
Schedulable Segments:	15				

- b. A non-scheduled only meeting environment does not require CTS-Manager, but does require users to manually dial or use speed dial entries to access meetings. When deploying a non-scheduled meeting environment for multipoint deployments with less than six rooms, it is recommended that a single speed dial entry be used for multipoint access (with five or less CTS endpoints it is only possible to have a single multipoint meeting). End users continue to use their existing calendaring system to reserve the rooms and use the single speed dial entry for multipoint meeting access.

Non-scheduled only meeting environment for TelePresence deployments with six or more rooms is not recommended. However, if this is the only option, the following consideration needs to be accounted for:

- Meeting access—How will users access meetings? With six or more TelePresence rooms it is possible to have multiple multipoint meetings. Providing multiple speed dials for multiple multipoint meetings may cause confusion and ultimately end up causing users to dial into the wrong meeting. In this environment users need to manually dial into multipoint meetings.

For this reason it is recommended that a centralized scheduler be used for multipoint meetings. The scheduler can either create a new static meeting for each request or select from a pool of pre-configured static meetings. After each request is processed by the scheduler, meeting information is sent to the meeting participants. At the time of the meeting, participants manually dial into the multipoint meeting using the information provided by the meeting scheduler.

Another option is to install a standalone directory and groupware server that is supported by the solution (e.g., Microsoft Active Directory and Exchange or IBM Domino). This allows a centralized scheduler to schedule multipoint meetings using Exchange which in conjunction with CTS-Manager provides system resource management for all scheduled multipoint meetings. This also provides one-button-to-push meeting access, eliminating any issues with users having to manually dial into meeting.

If ad hoc meetings are used, the scheduler launches the meeting at the scheduled start time, allowing users to walk into the room and attend their meeting. This is a very secure method for conducting multipoint meetings. CTS endpoints cannot dial into an ad hoc meeting; only the meeting administrator can add CTS endpoints to the meeting through the web GUI of the CTMS. However, this is very resource-intensive process considering every meeting must be manually initiated by the meeting administrator at the time of the meeting.

- Meeting security—Using a central scheduling resource to allocate static meetings is a resource-intensive process and prone to security risks. If a small number of multipoint numbers are used, it is possible for rouge endpoints to interrupt meetings. If a user tries to avoid the scheduling process and dials into the last multipoint meeting number they were assigned, they may interrupt a multipoint meeting in progress. To avoid this it is recommended that a maximum number of rooms be configured for each static meeting, minimizing the possibility of meeting interruptions. Configuring the number of rooms does not eliminate all potential risks. However it does minimize the threat by essentially locking the meeting after all the scheduled rooms are in the meeting. If a meeting requires more security, it is recommended that an ad hoc meeting be used. [Figure 11-11](#) illustrates an example of an eight room deployment with pre-configured meeting numbers.

Figure 11-11 Pre-configured Meeting Numbers

Deployment with 8 CTS endpoints



Number Range
54100 – 54300

54101 – 54170 = 3 room meeting
54171 – 54200 = 4 room meeting
54201 – 54250 = 5 room meeting
54251 – 54280 = 6 room meeting
54281 – 54290 = 7 room meeting
54291 – 54300 = 8 room meeting

224588

- Resource management—There is no resource management for non-scheduled meetings. In the event a centralized multipoint deployment supports more than 48 segments, a centralized scheduler is required to ensure multipoint resources are allocated properly. Maximum and ad hoc resources should be configured for the total number of available segments.
- c. Combined scheduled and non-scheduled meetings require CTS Manager. This deployment provides one-button-to-push meeting access for scheduled meetings and manual dial meeting access for non-scheduled meetings. This type of deployment allows personal static meeting numbers for power users or executives. These numbers can be used for last minute multipoint meetings when scheduling ahead of time is not convenient. Ad hoc meetings may also be used for high profile meetings or a white glove type meeting service. However, there are a number of considerations that must be taken into account:
- Meeting security—Since static meeting numbers are not secure, it is possible for an uninvited room to dial into the multipoint meeting.
 - Administrative resources—If static meetings are supported by a centralized scheduler, as described above, or ad hoc meetings are used, additional administrative resources are probably required.

- Resource management—There is no resource management for non-scheduled meetings. In the event a centralized multipoint deployment supports more than 48 segments, it is recommended that a separate CTMS devices be deployed, with one CTMS dedicated to scheduled meetings and one dedicated to non-scheduled meetings. This ensures that resources are always available for non-scheduled meetings.
3. Failover/redundancy—With the current release of CTMS and CTS Manager, automated failover is not supported. The following failover options are recommended for all meeting deployment scenarios described above:
 - a. Scheduled meeting deployment—In a scheduled meeting deployment, two CTMS devices can be configured in CTS Manager. CTMS-1 is configured in scheduled mode, while CTMS-2 is configured as non-scheduled. In case of a failure to CTMS-1, the system administrator uses the CTS Manger GUI interface to migrate all scheduled multipoint meetings to CTMS-2. The administrator then changes the control state of CTMS-1 to non-scheduled and changes the control state of CTMS-2 to scheduled. When meetings are migrated to CTMS-2, conference access number and meeting IDs are updated and new one button to push entries are propagated to all CTS endpoints. Any meeting that is in progress during the failure migration is not migrated.
 - b. Non-scheduled meeting deployment—For this deployment method a hot standby is recommended. Configure two CTMS devices with the same configuration, including the IP address. Manually shut down the Ethernet port to which CTMS-2 is connected. In case of a failure in CTMS-1, shut down the Ethernet port to which CTMS-1 is connected and **no shut** the Ethernet port to which CTMS-2 is connected.

Cisco Unified Communications Manager (CUCM) has the ability to route calls to a secondary CTMS in the event of a primary CTMS failure using route lists/route groups. However, this is not recommended, since there is no state information passed between CTMS devices. In the case of a temporary CTMS failure it is possible to have a meeting split between two CTMS devices (split meetings).

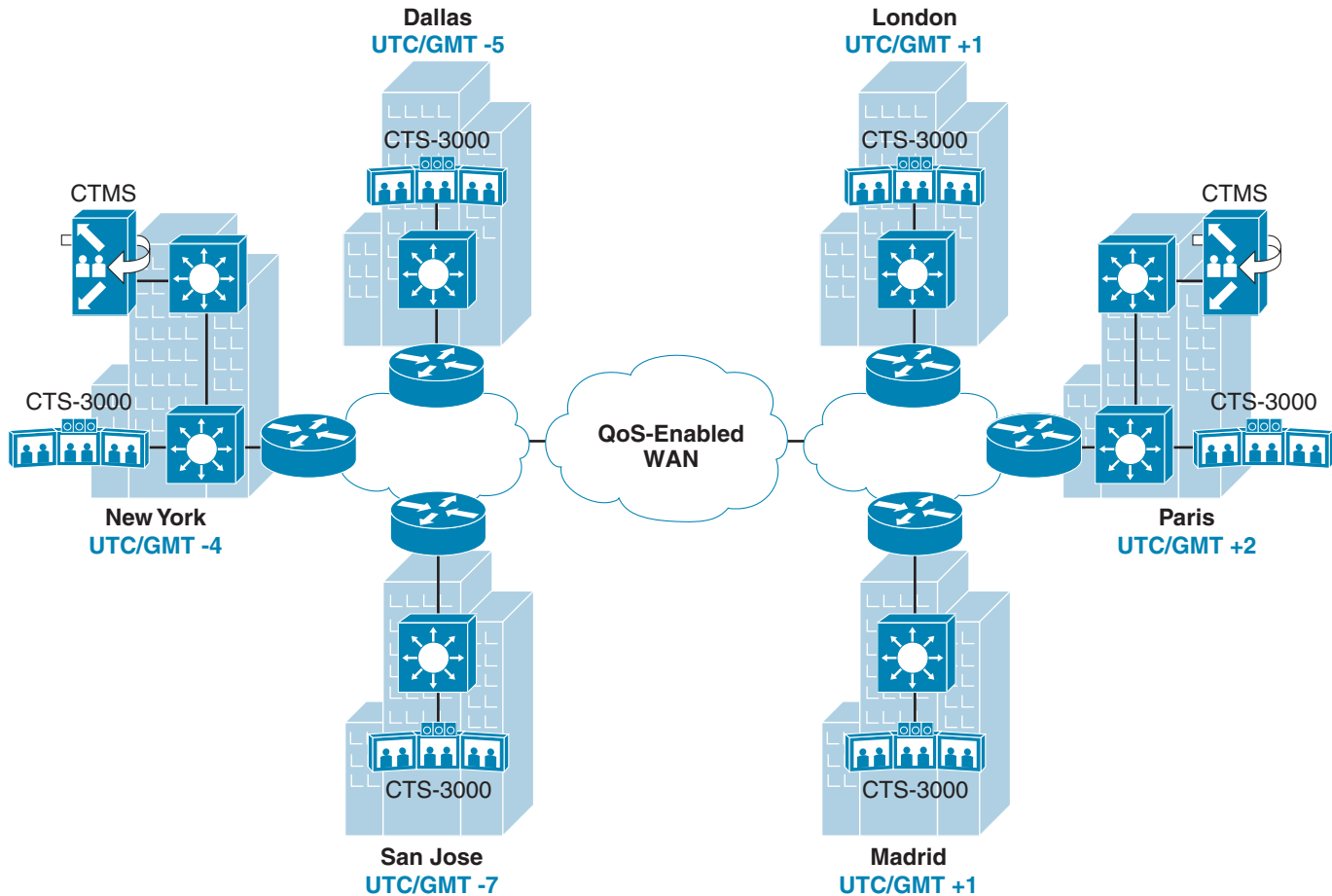
- c. Combined scheduled and non-scheduled—Almost all centralized deployments consist of less than 16 CTS-3000s, or a total of 48 table segments, allowing a single CTMS to accommodate all systems simultaneously. This does, however, provide a challenge for failover in an environment where a single CTMS is providing scheduled and non-scheduled resources. To provide seamless failover for all users, scheduled and non-scheduled resources must be supported on separate CTMS devices. Both failover methods described above for scheduled and non-scheduled must be used to supply failover.

Distributed Deployment

A distributed deployment is recommended for large TelePresence deployments or smaller deployments with three or more CTS endpoints in separate geographical regions, as shown in [Figure 11-12](#). As TelePresence networks grow, it is very advantageous to localize CTMS devices if possible.

Regionally localizing CTMS devices minimizes latency and saves bandwidth. [Figure 11-12](#) provides an example of a distributed deployment with a CTMS in New York providing multipoint services for North America and a CTMS in Paris providing multipoint services for Europe.

Figure 11-12 Distributed Multipoint Deployment



Note: UTC/GMT Times for North American Locations Based on U.S. Daylight Savings Time

224390



Note

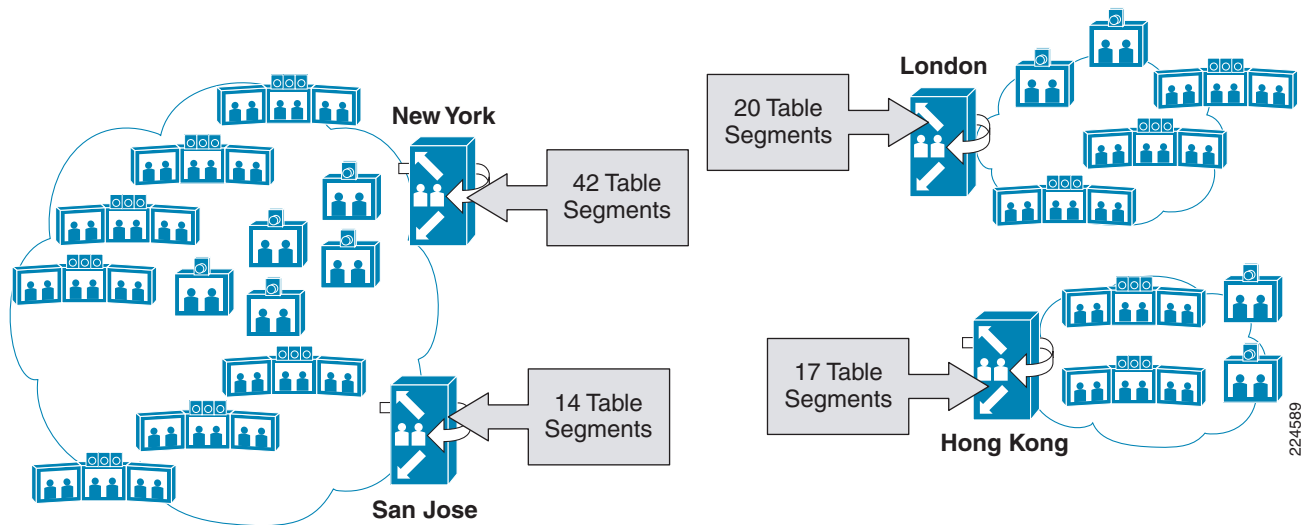
It should be noted that the current CTMS implementation (software version 1.1) does not support CTMS chaining/cascading for scalability.

Deployment Considerations

As seen above in the centralized deployment, there are a number of considerations that must be addressed for a simple multipoint deployment. In a distributed deployment, there are a number of additional considerations that must be addressed to ensure a successful deployment.

1. Selecting sites for the CTMS resources—As mentioned above, multipoint switches need to be located at a sites providing end-to-end network latency less than 200ms, for targeted CTS endpoints, and adequate bandwidth for the number of CTS segments supported by each site. In [Figure 11-13](#), four regional locations have been selected for CTMS resources based on the number of regional CTS endpoints, bandwidth availability, and proximity to the regional CTS endpoints.

Figure 11-13 Distributed CTMS Deployment Example



2. Supported Meeting Types—In a distributed multipoint deployment, it is recommended that scheduled only or a combination scheduled and non-scheduled meeting environment be supported. Non-scheduled only meeting environments should be avoided in distributed multipoint deployments due to the complexity of administering multipoint meetings.

- a. A scheduled only meeting environment requires CTS Manager and provides one-button-to-push meeting access for end users. Meetings are scheduled by end users or a centralized scheduling group using Microsoft Exchange or IBM Domino. When configuring CTMS resources, maximum segments should be configured for the total available segments plus additional ad hoc segments. It is important to remember that ad hoc resources need to be available even in a scheduled only environment. If no ad hoc resources are configured on the CTMS, there is no way to add non-scheduled CTS endpoints to scheduled meetings (as described in [Multipoint Resources](#) in [Chapter 10, “Cisco TelePresence Multipoint Solution Essentials”](#)).
- b. Non-scheduled meeting environments are not recommended in distributed multipoint deployments. Manually managing resources, meeting placement, and scheduling meetings is a difficult task that is prone to errors.

CTS Manager is required to support scheduled meetings and manage CTMS resources. If a calendaring system other than Exchange or Domino is used to schedule meetings/rooms, it is recommended that a standalone Active Directory and Exchange server be deployed. This allows a centralized scheduler to schedule multipoint meetings using Exchange which, in conjunction with CTS Manager, provides system and geographical resource management for all scheduled multipoint meetings. This also provides one-button-to-push meeting access, eliminating any issue with users having to manually dial into meetings.

- c. Combined scheduled and non-scheduled meetings require CTS Manager. This deployment provides one button to push meeting access for scheduled meetings and manual dial meeting access for non-scheduled meetings. This type of deployment allows personal static meeting numbers for power users or executives. These numbers can be used for last minute multipoint meetings when scheduling ahead of time is not convenient. Ad hoc meetings may also be used for high profile meetings or a white glove type meeting service. However, there are a number of considerations that must be taken into account:
 - Meeting security—Since static meeting numbers are not secure, it is possible for an uninvited room to dial into the multipoint meeting.

- Administrative resources—If static meetings are supported by a centralized scheduler, as described [Centralized Deployment](#), or ad hoc meetings are used, additional administrative resources are probably required.
 - Resource management—There is no resource management for non-scheduled meetings. For this reason, it is recommended that scheduled and non-scheduled resources be supported on separate CTMS devices. This ensures that enough resources are available for non-scheduled meetings and executives do not call complaining about their personal multipoint number not working.
3. CTMS resources per site—After determining which sites will provide multipoint resources, it is important to determine how many segments each site will support. [Table 11-1](#) shows the segment breakdown for each multipoint site, based on [Figure 11-13](#).

Table 11-1 Resource Allocation

Site	Total Segments	Scheduled Segments	Ad hoc Segments
New York	42	36	6
San Jose	17	14	3
London	20	17	3
Hong Kong	17	14	3

Resources for each site are broken down into three categories, total segments, scheduled segments, and ad hoc segments. Total segments is used to limit the number of connections each CTMS device supports. This ensures that the bandwidth allocated to a multipoint site is not exceeded. Scheduled segments is passed to CTS Manager and used to manage resources on each CTMS device. Ad hoc Segments is used to add non-scheduled CTS endpoints to scheduled meetings or to provide resources for static and ad hoc meetings. In [Table 11-1](#), ad hoc segments are only being provided to add non-scheduled CTS endpoints to scheduled meetings.

Determining how many multipoint resources are assigned to each site is based on call patterns and WAN bandwidth. There is no exact calculation for determining the number of multipoint resources for a region. However, at a minimum there should be enough resources allocated to support all CTS endpoints within the region.

In [Table 11-1](#), it is decided that New York is the hub site, providing enough schedulable resources for a single multipoint meeting containing all deployed CTS-3000s. Six ad hoc resources are added, allowing non-scheduled CTS endpoints to be added to scheduled meetings. Total segments is configured for 42 to ensure the provisioned bandwidth for the hub site is not exceeded.

San Jose is configured with enough scheduled resources to support all regional CTS endpoints and six additional CTS segments. Three ad hoc resources are added, allowing non-scheduled CTS endpoints to be added to scheduled meetings. Total segments is configured for 17 to ensure the provisioned bandwidth for Hong Kong is not exceeded.

London is configured with enough scheduled resources to support all regional CTS endpoints and six additional CTS segments. Three ad hoc resources are added, allowing non-scheduled CTS endpoints to be added to scheduled meetings. Total Ssegments is configured for 20 to ensure the provisioned bandwidth for London is not exceeded.

4. Required bandwidth for each multipoint site:

The appropriate bandwidth must be configured for each multipoint site. Calculate the required bandwidth for each multipoint site using the calculations from the previous section outlining bandwidth requirements.

5. CTMS configurations for geographical selection:

As described in [Geographical Resource Management](#) in [Chapter 10, “Cisco TelePresence Multipoint Solution Essentials”](#), CTS Manager provides the ability to select a CTMS device with the closest proximity to scheduled CTS endpoints. It is important to carefully analyze the location of CTS endpoints and the regional multipoint sites to determine the best time zone entry for each CTMS. CTMS time zone entries may have to be modified to obtain the most accurate meeting placement.

6. Failover/redundancy—With the current release of CTMS and CTS Manager, automated failover is not supported. The following failover options are recommended for all meeting deployment scenarios described above:

- a. Scheduled meeting deployment—In a scheduled meeting deployment, two CTMS devices can be configured in CTS Manager. CTMS-1 is configured in scheduled mode, while CTMS-2 is configured as non-scheduled. In case of a failure to CTMS-1, the system administrator uses the CTS Manager GUI interface to migrate all scheduled multipoint meetings to CTMS-2. The administrator then changes the control state of CTMS-1 to non-scheduled and changes the control state of CTMS-2 to scheduled. When meetings are migrated to CTMS-2, conference access number and meeting IDs are updated and new one button to push entries are propagated to all CTS endpoints. Any meeting that is in progress during the failure migration is not migrated.
- b. Non-scheduled meeting deployment—Non-scheduled only meeting deployments are not recommended in a distributed deployment.
- c. Combined scheduled and non-scheduled—It is recommended that scheduled and non-scheduled multipoint resources be supported on separate CTMS devices. Failover scenarios for each device should be the same as described above for scheduled and non-scheduled meetings.

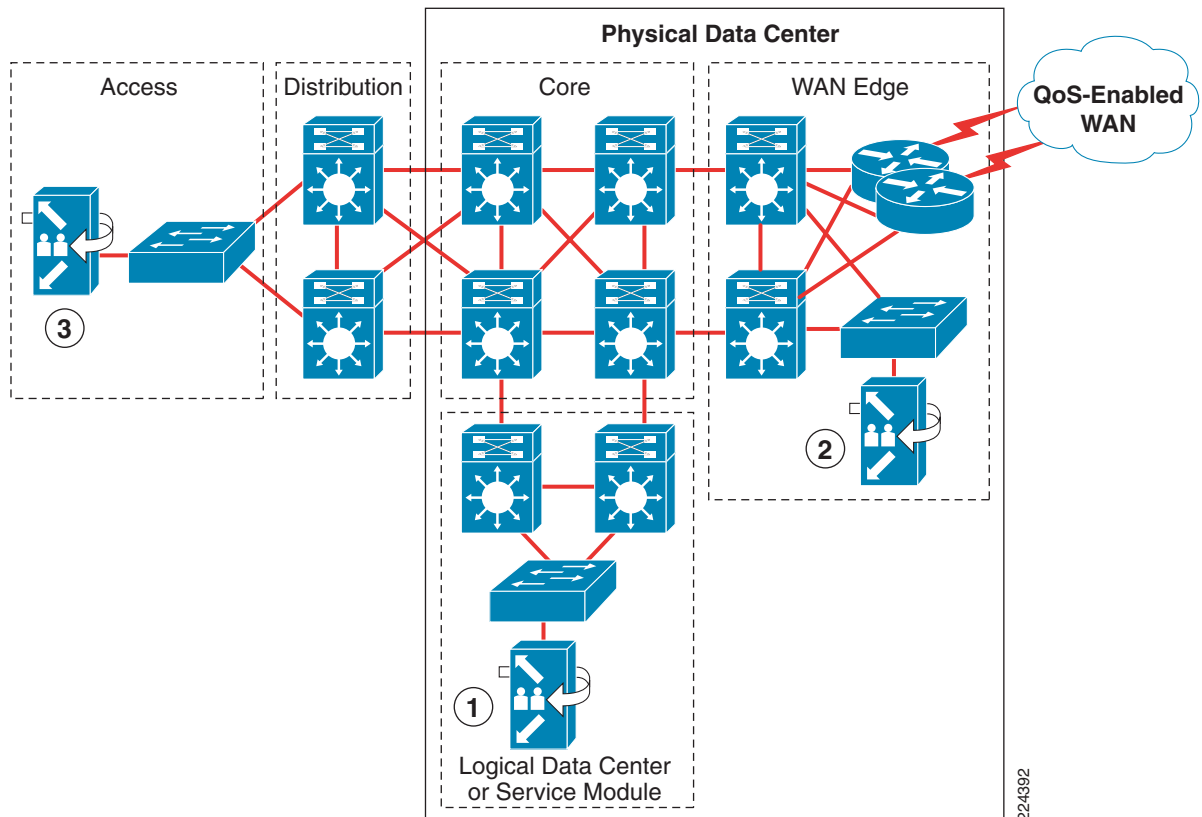
Positioning of the CTMS within the Campus or Branch

Due to the total audio and video bandwidth requirements for a Cisco TelePresence CTMS (which can be up to approximately 260 Mbps), it is important to consider its placement in the network. Within a campus deployment, placement of the CTMS within a logical data center LAN segment may be desirable due to the availability of bandwidth, an uninterruptible power supply, as well as ease of monitoring. The downside is that all multipoint TelePresence traffic must be backhauled into and out of the logical data center LAN segment. The data center design may need to be adjusted to accommodate the necessary increase in traffic.

An alternative is to locate the CTMS at the access layer, towards the logical WAN edge of the campus. Note that some customers may not have an access layer switch at this location. This type of placement may minimize the amount of traffic which is backhauled through the campus LAN, however this is dependent upon the location of the CTS endpoints. If the majority of the CTS endpoints within the multipoint TelePresence deployment are remote to the campus location, this design may provide some benefit. If the majority of the CTS endpoints are within the campus, this design may provide little benefit. The downside to this placement is that an uninterruptible power supply may not be available, depending upon the physical location of WAN network devices within the campus deployment. However, in many campus network deployments, the WAN routers are physically located within a data center itself, although not logically on a data center LAN segment.

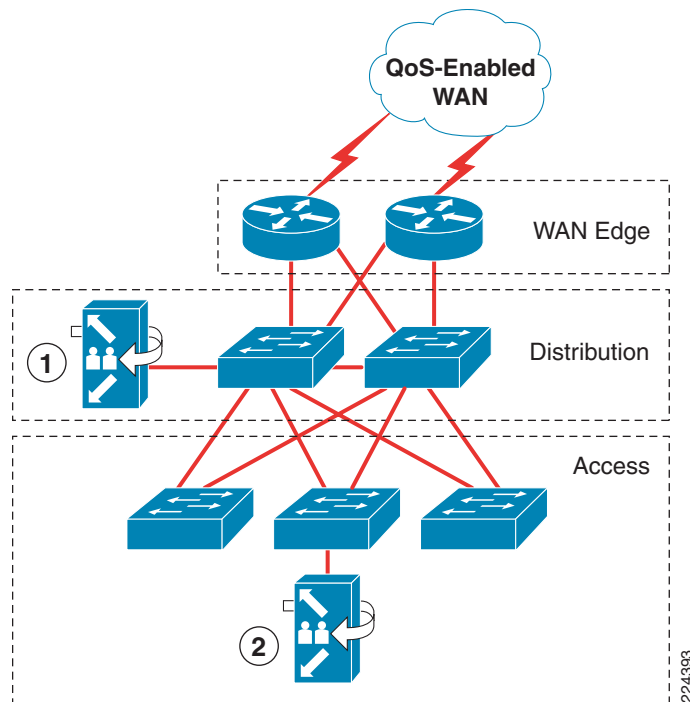
A third alternative is to simply locate the CTMS at the access layer within the campus network. This type of placement minimizes the amount of unnecessary traffic to the logical WAN edge and the logical data center if the majority of CTS endpoints are within the campus. The downside is that the likelihood of an uninterruptible power for the CTMS may be lower at the access layer. Figure 11-14 shows the three campus placement alternatives.

Figure 11-14 Possible CTMS Locations within the Campus



Under some circumstances, it may be necessary to deploy a CTMS at a branch location. However, due to limited bandwidth of branch locations, this design is not highly recommended. When deploying at a branch, it is recommended that the CTMS be deployed at the distribution layer of any hierarchical LAN configuration, as shown in Figure 11-15.

Figure 11-15 Possible CTMS Locations Within the Branch



This minimizes the amount of unnecessary traffic backhauled through the branch LAN network. An alternative is to place the CTMS at the access layer if no available LAN ports exist at the distribution layer.

Next, it should be noted that the CTMS currently supports a single 1 Gigabit Ethernet connection. Resilient connections to dual LAN switches are currently not supported. Further, since the CTMS is capable of generating traffic loads in excess of 100 Mbps, it is not recommended to place the CTMS on a 100 Mbps Ethernet LAN port.

Finally, both CTS endpoints and their associated IP phones transmit CDP packets to associated network devices (switches and routers). CDP can be used to extend trust to the device connected to the LAN switch port and automatically place CTS endpoints within a voice VLAN. If the CTS endpoints within a network deployment are located within a separate voice VLAN, placement of the CTMS within the voice VLAN maintains consistency of the overall TelePresence deployment from a traffic isolation and QoS viewpoint.

Network Requirements

Maintaining the required SLAs for Cisco TelePresence can be challenging when multipoint is added to the network. SLA requirements defined for point-to-point TelePresence in [Chapter 4, “Quality of Service Design for TelePresence”](#) and duplicated in [Table 11-2](#) remain the same and should not be affected with the addition of multipoint. However, providing acceptable latency for multipoint meetings can be a challenge for geographically disperse deployments. The three main network considerations that need to be carefully considered when deploying multipoint capabilities for Cisco TelePresence are bandwidth, latency, and traffic bursts. Latency and bandwidth are discussed in this section. [Estimating Burst Sizes within Multipoint TelePresence Calls](#) thoroughly discusses bursts within a multipoint TelePresence design. How and where the CTMS is deployed on the network directly affects latency for

multipoint meetings and bandwidth patterns on the network. Deploying a multipoint device in the wrong location, physical or geographical, may cause an undesirable meeting experience and directly affect network performance.

Table 11-2 Cisco TelePresence SLA Requirements

Metric	Target	Threshold 1 (Warning)	Threshold 2 (Call Drop)	Enterprise Component	Service Provider Component
Latency	150 ms	200ms	400 ms ¹	20%	80%
Jitter	10 ms	20 ms	40 ms	50%	50%
Loss	0.05%	0.10%	0.20%	50%	50%

1. Call is not dropped.

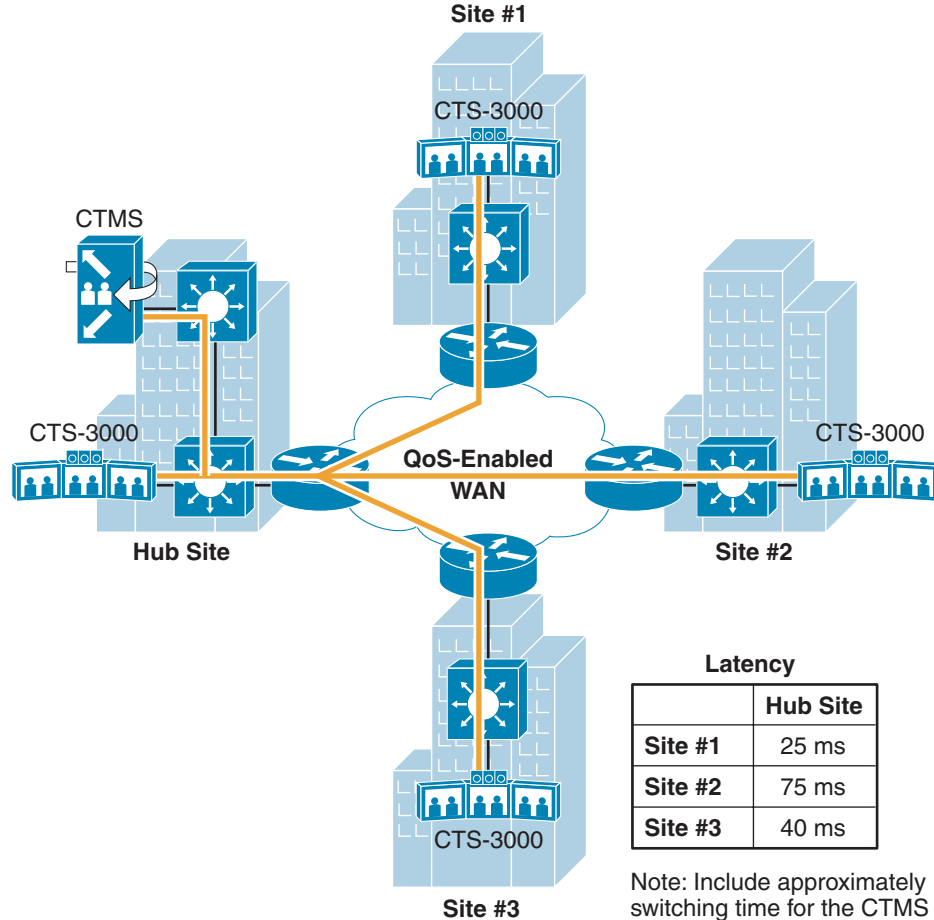
Latency

One of the key differentiators for Cisco TelePresence is the ability to maintain extremely low latency while providing high quality 1080p video and spatial audio. Excessive latency in any Cisco TelePresence meeting will degrade the “in-person” experience. Latency becomes an even bigger issue with multipoint, since all CTS systems dial into a CTMS that may not be located in the same geographic location as the CTS endpoints. Due to the nature of multipoint, two CTS endpoints that provide very low latency in a point-to-point meeting may have considerably higher latency in a multipoint meeting. Inserting any multipoint device in the media path of a Cisco TelePresence call introduces additional latency. However, proper placement of the CTMS helps minimize latency and preserve the Cisco TelePresence experience.

A Cisco TelePresence network should always be designed to target one-way, end-to-end, network latency of less than <150ms. However, in some cases this is not possible due to long distances between international sites. Therefore, the upper limit allowed for one-way, end-to-end network latency is < 200ms. Anything above 200ms causes the message “Experiencing Network Delay” to be displayed on the Cisco TelePresence endpoint, degrading the user experience. Therefore, multipoint deployments should provide one-way, end-to-end, latency below 200ms in all cases.

Figure 11-16 illustrates a three site multipoint deployment with the CTMS located in the hub site.

Figure 11-16 Multipoint Network Latency Example



224394

As mentioned above, one way, network only latency must stay below 200ms to maintain the TelePresence experience. In Figure 11-16, the hub site is chosen to deploy the CTMS. Using the latency matrix in the diagram, the highest latency for any multipoint meeting is between Site #2 and Site #3. To calculate the “highest” latency for a multipoint deployment, take the two sites with the highest latency between themselves and the CTMS and add 10ms for CTMS switching delay. The worst case latency in Figure 11-16 is calculated as:

$$\text{Site \#2 - Hub } 75\text{ms} + \text{Site\#3 - Hub } 40\text{ms} + \text{CTMS } 10\text{ms} = 125\text{ms}$$

As previously noted, the current CTMS implementation (software version 1.1) does not support chaining/cascading for scalability. Therefore, the latency examples above also apply to distributed multipoint designs. The design engineer should also be aware that any CTS endpoint could potentially utilize any CTMS within the network as part of a multipoint call and should be aware of the maximum possible end-to-end latency of any combination of CTS endpoints utilizing any CTMS within the network.

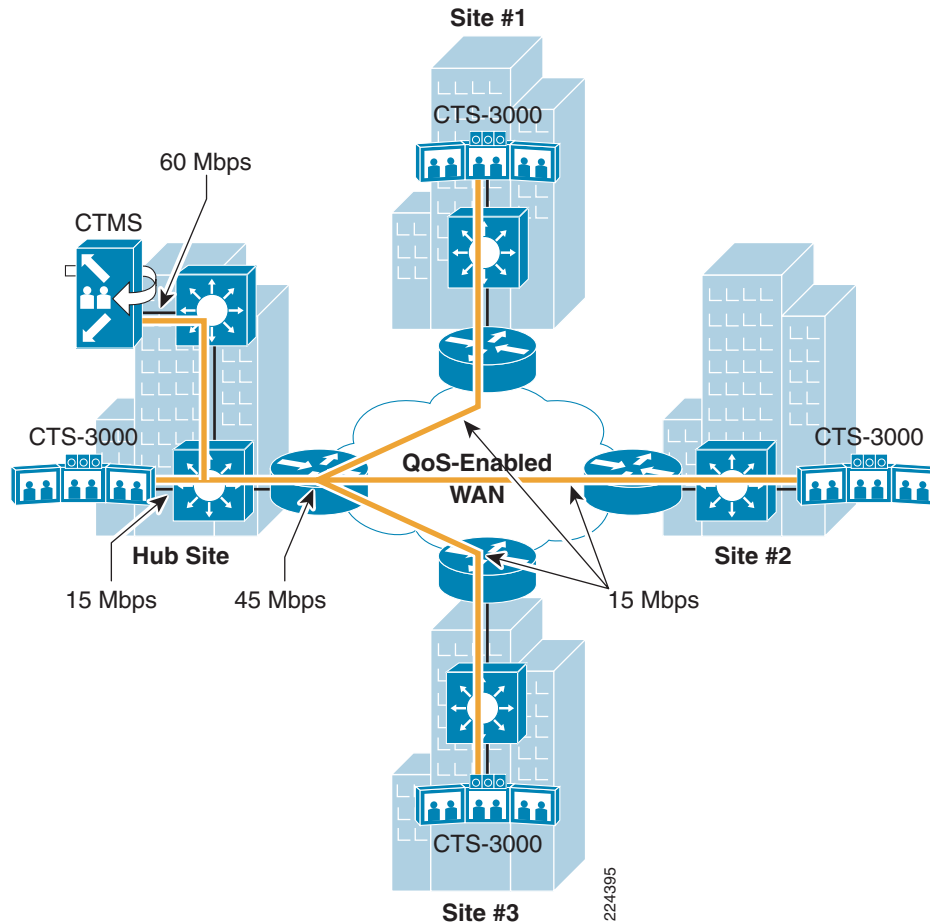
Bandwidth

Sufficient bandwidth must also be provisioned to the site which houses the CTMS to support the additional traffic required for multipoint meetings, as well as the traffic required for point-to-point meetings.

Centralized Multipoint Designs

Figure 11-17 highlights the bandwidth requirements of centralized multipoint designs by extending the simple multisite example shown in Figure 11-9 to include the bandwidth requirements for TelePresence for each site.

Figure 11-17 Centralized Multipoint Design Bandwidth Example



Calculating the amount of bandwidth required at the Hub Site is fairly straight forward. In the example above, the circuit to the Hub Site must have sufficient bandwidth to support three CTS-3000 systems at 1080p even though the Hub Site only houses a single CTS-3000. This is due to the CTMS being located at the Hub Site. Audio and video traffic from Site #1, Site #2, and Site #3 must traverse the circuit during multipoint meetings. Note that the LAN infrastructure within the Hub Site must also be designed to support the cumulative bandwidth of all four TelePresence CTS endpoints.

A general rule of thumb for 1080p configurations is to simply estimate 15 Mbps per CTS-3000 or CTS-3200 and 5.5 Mbps per CTS-1000 or CTS-500 with low-speed auxiliary video input.



Note

[Audio and Video Flows In A Multipoint TelePresence Design](#) presents a thorough discussion on how to calculate the audio and video flows to and from a CTMS in a multipoint meeting; since traffic flows are asymmetric in a multipoint call.

Provisioning the correct amount of bandwidth on the LAN and WAN is essential for a successful multipoint deployment. As illustrated in Figure 11-17, the maximum potential bandwidth for each CTS-3000 (15Mbps) is provisioned to ensure there is no packet loss due to insufficient bandwidth. However, the actual bandwidth used during a multipoint meeting typically averages less (10 - 12 Mbps with six people sitting at the table participating in the meeting) than the provisioned bandwidth.

The design engineer should also keep in mind that the amount of bandwidth provisioned to the site which houses the CTMS must be increased for point-to-point meetings which occur at the same time as multipoint meetings.

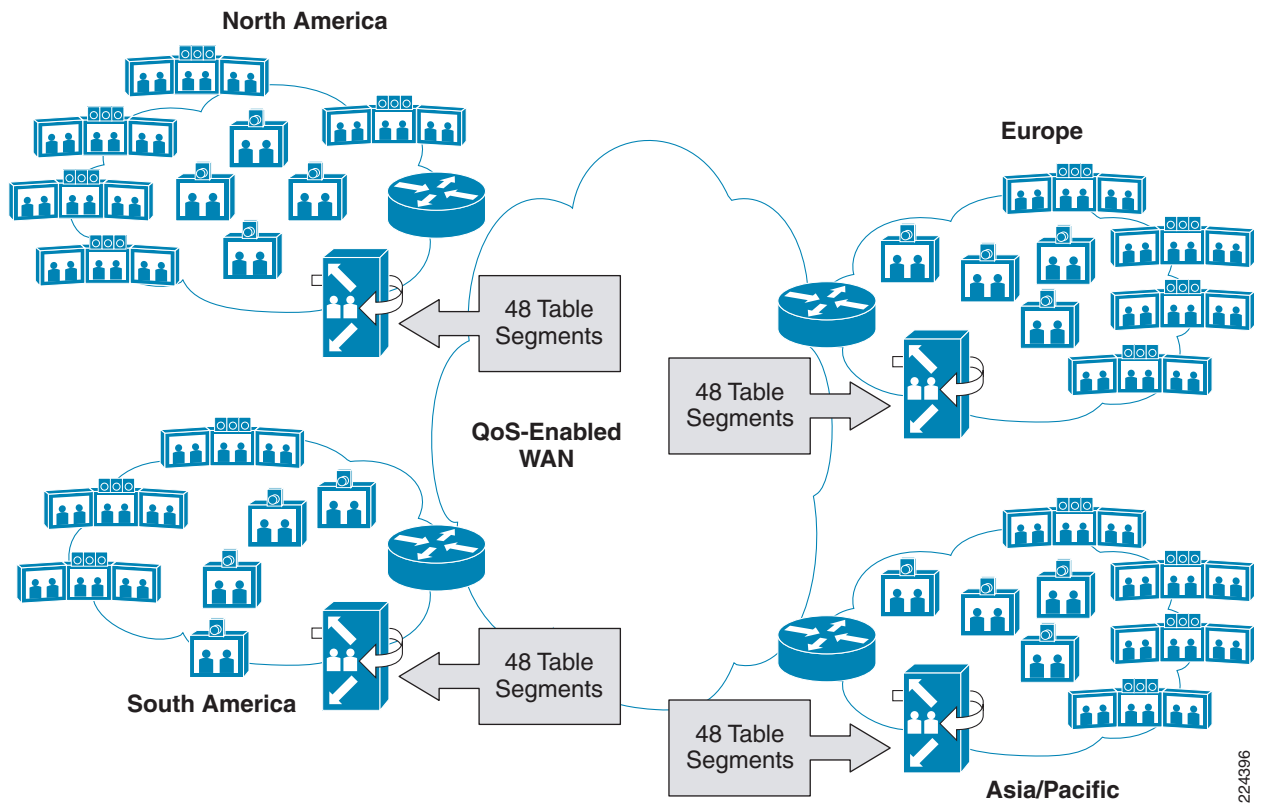
Distributed Multipoint Designs

As illustrated above, provisioning bandwidth for TelePresence deployments with a single CTMS and a limited number of CTS systems is fairly straight forward. However, in larger deployments with multiple CTMS devices and a mix of CTS3000s, CTS-3200s, CTS-1000s, and CTS500s, bandwidth provisioning becomes more difficult. Several methods of provisioning bandwidth at the CTMS site are explored within the next sections.

Maximum Bandwidth Per CTMS Approach

Figure 11-18 provides an example of a large TelePresence network with distributed CTMS devices. The CTMS devices are deployed in regions around the world. In this example, the network design engineer must determine how much bandwidth needs to be provisioned to each CTMS site to handle not only multipoint calls involving CTS endpoints within the region, but also CTS endpoints which may be across the QoS-enabled WAN as well.

Figure 11-18 Bandwidth Provisioning Based Maximum Configured CTMS Capacity



The approach shown in [Figure 11-18](#) is to simply provision sufficient bandwidth to accommodate the maximum amount of traffic from the CTMS at each regional location, assuming each CTMS is configured for its maximum of 48 segments (as of software version 1.1) or configured for less than the maximum segments. This method may be beneficial to customer deployments in which the number of CTS table segments (1 segment per CTS-1000 or CTS-500 and 3 segments per CTS-3000 or CTS-3200) deployed throughout the network greatly exceeds the maximum capacity of a single CTMS. In this type of large deployment, the customer may have no control of which sites have multipoint meetings with each other and what type of CTS units are involved in the meeting.

A rough estimate for calculating the required bandwidth is to simply multiply 5.5 Mbps per CTS-1000 by the maximum number of segments (also referred to as table segments) supported by the CTMS:

$$5.5 \text{ Mbps} \times 48 \text{ table segments} = 264 \text{ Mbps}$$

For example, if every CTS unit in the North America, South America, and Europe regions shown in [Figure 11-18](#) were configured for a single multipoint call, the CTMS unit selected for the call and the bandwidth provisioned for TelePresence to the regional site which housed the CTMS would need to handle 48 segments.

Keep in mind that additional bandwidth capacity is required to handle additional point-to-point calls to and from the regional site as well if there were more CTS units at the regional site not involved in the multipoint call.

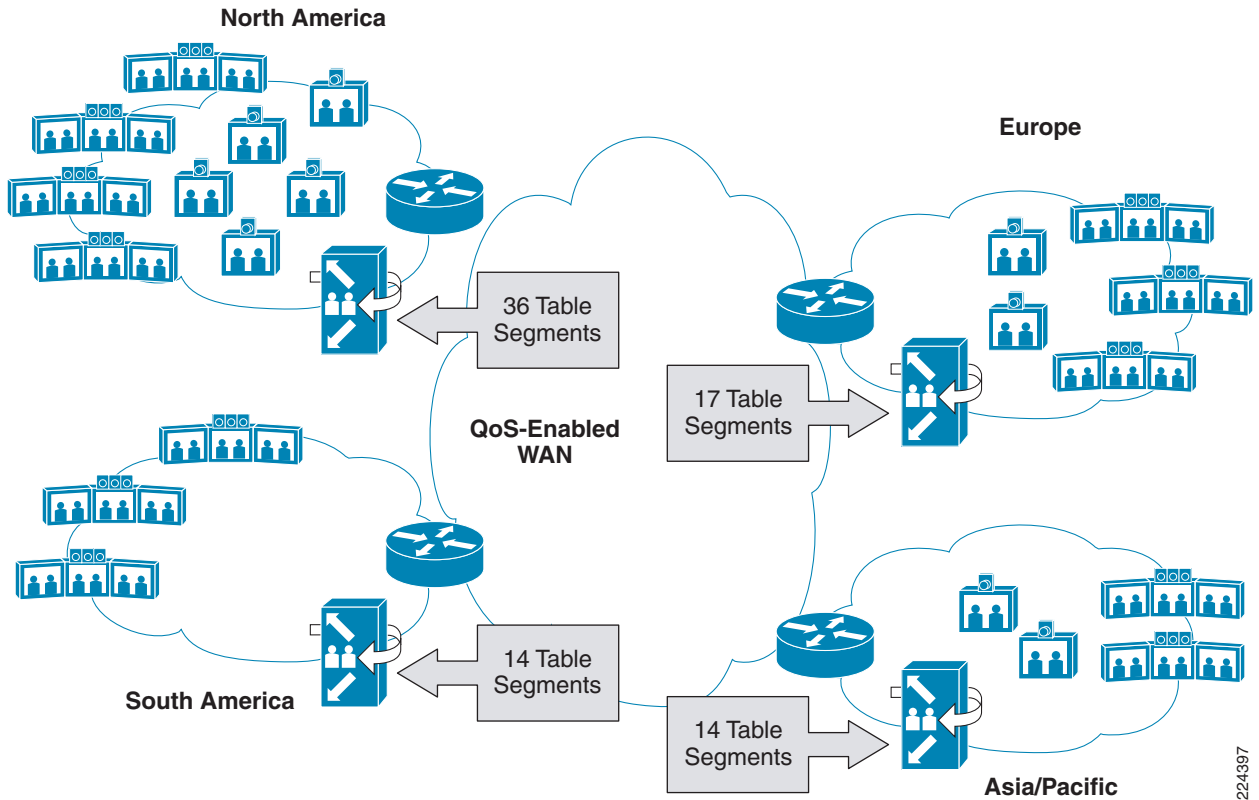
The advantage of this method of bandwidth provisioning is that as networks grow, the network design engineer does not constantly have to increase bandwidth to each site. The downside is that for many customers, provisioning that much bandwidth is unfeasible from a cost perspective.

Bandwidth Allocation Based on Meeting Patterns

A second method of provisioning bandwidth is based on historical meeting patterns and knowledge of the specific CTS units within the network. This method relies on limiting the Maximum Segments defined within the CTMS to be at or below the bandwidth allocated for TelePresence meetings from that regional location. This method may be beneficial to customer deployments in which the total number of CTS table segments may not exceed the maximum capacity of a single CTMS.

An example of this type of distributed multipoint TelePresence design is shown in [Figure 11-19](#).

Figure 11-19 Bandwidth Provisioning Based Historical Meeting Patterns



In the example shown in [Figure 11-19](#), based on WAN bandwidth and meeting patterns, the North America CTMS device is configured with a maximum of 36 table segments. Based on WAN bandwidth and meeting patterns, the Europe CTMS device is configured with the maximum 17 table segments. The remaining multipoint devices are configured to support a limited number of table segments.

When provisioning bandwidth using this method, it may be beneficial to base the bandwidth calculations for each site on the type of CTS system and video resolution supported, in order to more accurately assess the bandwidth requirement. A CTS-1000 requires more bandwidth than a CTS-3000 when provisioning is based on table segment versus entire system. A rough bandwidth guideline is to allocate 5.5 Mbps for a CTS-1000 system (one table segment) and 15 Mbps (three table segments at 5 Mbps each) for a CTS-3000 system. [Table 11-3](#) provides the bandwidth required per system for each resolution and motion handling.

Table 11-3 Bandwidth Provisioning

Resolution	1080p	1080p	1080p	720p	720p	720p
Motion Handling	Best	Better	Good	Best	Better	Good
CTS-1000 Provisioned Bandwidth for Multipoint	5.5 Mbps	4.9 Mbps	4.4 Mbps	4.4 Mbps	3.2 Mbps	2.1 Mbps
CTS-3000 Provisioned Bandwidth for Multipoint	15 Mbps	12.9 Mbps	11.3 Mbps	11.3 Mbps	7.8 Mbps	4.4 Mbps

Below is the calculation for each multipoint device in [Figure 11-19](#). The calculations are based on all CTS endpoints running at 1080p best resolution. In environments with mixed resolutions, it is recommended that bandwidth be provisioned based on the highest resolution only.

North America region:

9 - CTS-1000 @ 5.5Mbps = 49.5 Mbps	9 table segments
<u>9 - CTS-3000 @ 15Mbps = 135 Mbps</u>	<u>27 table segments</u>
Multipoint bandwidth = 184.5 Mbps	36 table segments

South America and Asia/Pacific regions:

8 - CTS-1000 5.5Mbps = 44 Mbps	8 table segments
<u>2 - CTS-3000 15Mbps = 30 Mbps</u>	<u>6 table segments</u>
Multipoint bandwidth = 74 Mbps	14 table segments

Europe region:

8 - CTS-1000 5.5Mbps = 44 Mbps	8 table segments
<u>3 - CTS-3000 15Mbps = 45 Mbps</u>	<u>9 table segments</u>
Multipoint bandwidth = 89 Mbps	17 table segments



Note

The calculations above are only for multipoint calls. Additional WAN bandwidth must be provisioned for CTS systems participating in point-to-point calls located at sites with CTMS devices.

Estimating Burst Sizes within Multipoint TelePresence Calls

This section presents a brief discussion of the causes of bursts within a multipoint TelePresence call; and a method of estimating the size of those bursts so that the network can be provisioned accordingly.

Causes of Bursts within Multipoint TelePresence Calls

In multipoint TelePresence calls, bursts are generated as a result of one of the following events.

- Whenever a CTS endpoint joins a multipoint call
 - If the call is configured with the Video Announce feature enabled, when a new CTS endpoint joins the call, it becomes the active site. This causes an I-frame to be generated by the new CTS endpoint and replicated to every other CTS endpoint by the CTMS. In addition, the last active site sends an I-frame to the CTMS which is replicated and sent to the new CTS endpoint. Note that for CTS endpoints with multiple screens, multiple I-frames may be generated.
 - If the call is configured with the Video Announce feature disabled, it does not become the active speaker when joining the call. However, it needs to receive an I-frame to begin displaying video. Therefore, the active site sends an I-frame which is replicated by the CTMS to every CTS endpoint in the call. Note again that for CTS endpoints with multiple screens, multiple I-frames may still be generated.
- Normal transitioning of the active site or segment from one CTS endpoint to another CTS endpoint
 - This causes an I-frame to be generated by the new active site or segment and replicated by the CTMS to all of the other CTS endpoints in the multipoint call.

- One or more CTS endpoints or codecs reports loss in the received video
This causes the active site or segment to generate an I-frame to resynchronize all CTS endpoints. The I-frame is replicated by the CTMS and sent to all CTS endpoints in the multipoint call.
- Periodic synchronization of the CTS endpoints by the active site
Each active site or segment periodically sends out a new I-frame to synchronize the CTS endpoints. This occurs approximately every 5 minutes with the current TelePresence solution. The I-frame is replicated by the CTMS and sent to all CTS endpoints in the multipoint call.
- Normal transitioning of the video input from a device connected to the auxiliary input of one of the CTS endpoints.
Only one CTS endpoint at a time can function as a “presenter” within a multipoint meeting through the use of the Auto-Collaborate feature. Whenever the presenting CTS endpoint changes the content of the auxiliary video input, such as transitioning a PowerPoint slide, the burst of content is replicated by the CTMS to all of the other CTS endpoints within the multipoint call.
- Normal replication of the P-frame video by the CTMS
The active site or segment sends a new video frame (P-frame) every 33 ms when not sending an I-frame. This is replicated by the CTMS and sent to all of the other CTS endpoints within the multipoint call. Likewise, the last active sites or segments send P-frames every 33 ms.

Bursts Due to I-Frame Replication

The size of I-frames generated during speaker transitions, periodic synchronization of the video, synchronization of the video due to packet loss, and when new CTS endpoints join the multipoint call is highly variable. However, under normal lighting and background conditions, ESE testing has shown that the maximum size of I-frames generated by these events is approximately 64 Kbytes.

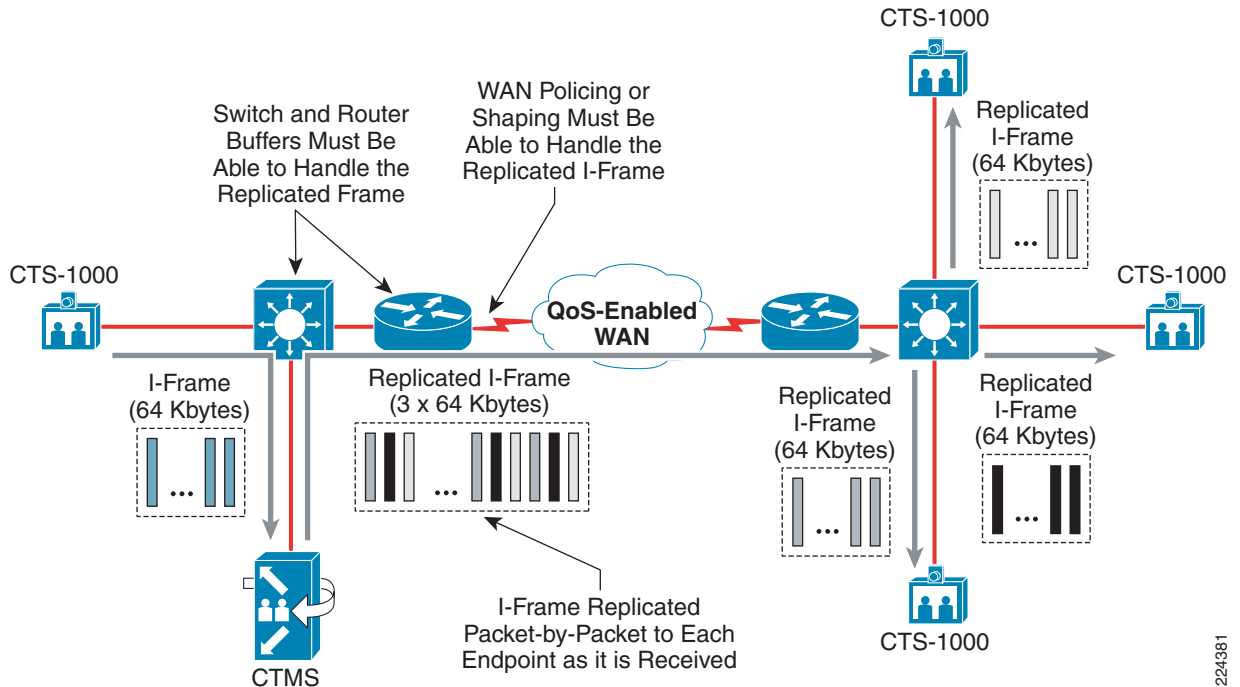


Note

ESE testing has also indicated that burst sizes may exceed 64 Kbytes if lighting and background conditions do not comply with Cisco recommendations. It is therefore critical to follow Cisco documented room design recommendations in order to ensure proper functioning of the TelePresence deployment over the network.

This burst is replicated on a packet-by-packet basis by the CTMS to every endpoint in the multipoint call, as shown in [Figure 11-20](#).

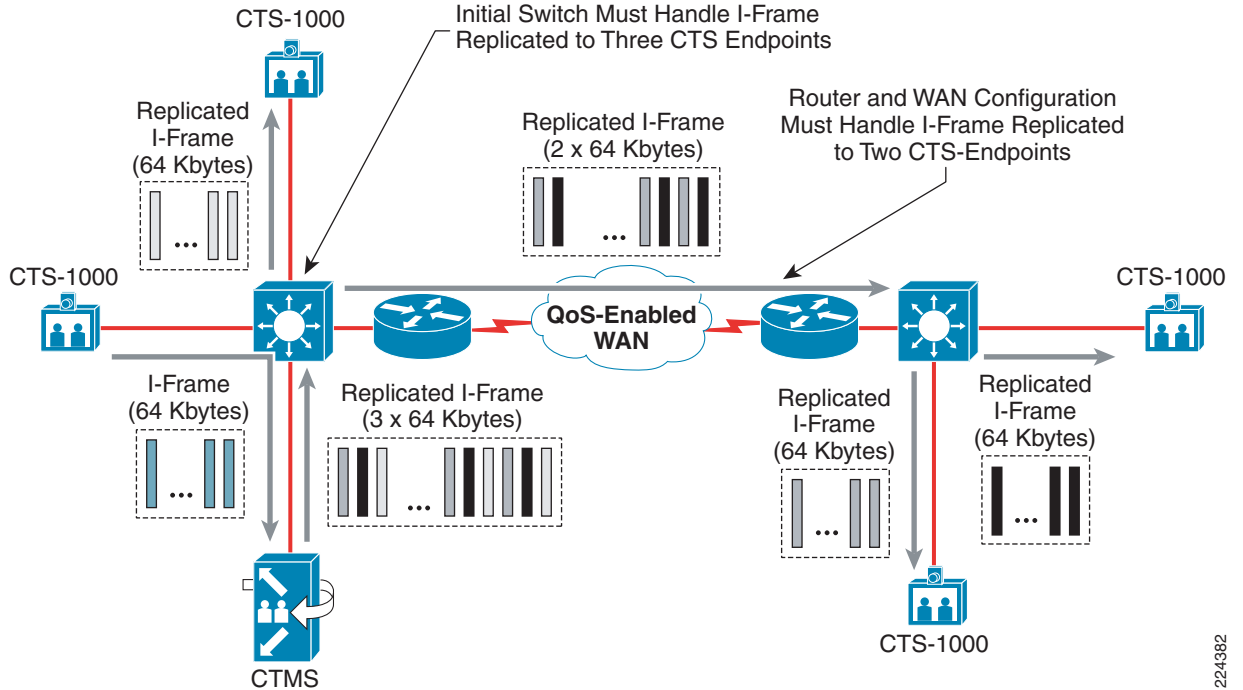
Figure 11-20 Burst Size Due to I-Frame Replication by the CTMS



It can be seen that as the number of CTS endpoints in the call increases, the size of the replicated I-frame burst also increases. In the example above the size of the I-frame burst sent by the CTMS is approximately $3 \times 64 \text{ Kbytes} = 192 \text{ Kbytes}$. This is sent within a single frame interval of 33 ms. Switch and router buffers need to be large enough to accommodate these bursts. Also, burst parameters configured within policers and/or shapers on any WAN circuits must also be large enough to accommodate the replicated I-frames.

Location of the CTS Endpoints

The location of the CTS endpoints within the multipoint call influences the size of the replicated I-frame burst across any given network device or WAN circuit. Figure 11-21 shows the same four-site multipoint call, with one of the CTS endpoints moved back to the head-end location.

Figure 11-21 Burst Size Due to Location of CTS Endpoints

As can be seen from [Figure 11-21](#), the switch to which the CTMS is connected must still have sufficient buffer capacity on the ingress port to handle the I-frame replicated to three CTS endpoints. However, only two CTS-endpoints are now located across the WAN. The size of the I-frame burst sent across the WAN is now approximately $2 \times 64 \text{ Kbytes} = 128 \text{ Kbytes}$. Again, this is sent within a single frame interval of 33 ms. Therefore, the egress port of the switch which serves as the uplink between the switch and the router, as well as the router LAN ingress port, must have sufficient buffer capacity to handle the I-frame replicated to two CTS endpoints. Any policer and/or shaper parameters configured on the router or within the WAN must now be able to accommodate the I-frame replicated to two CTS endpoints.

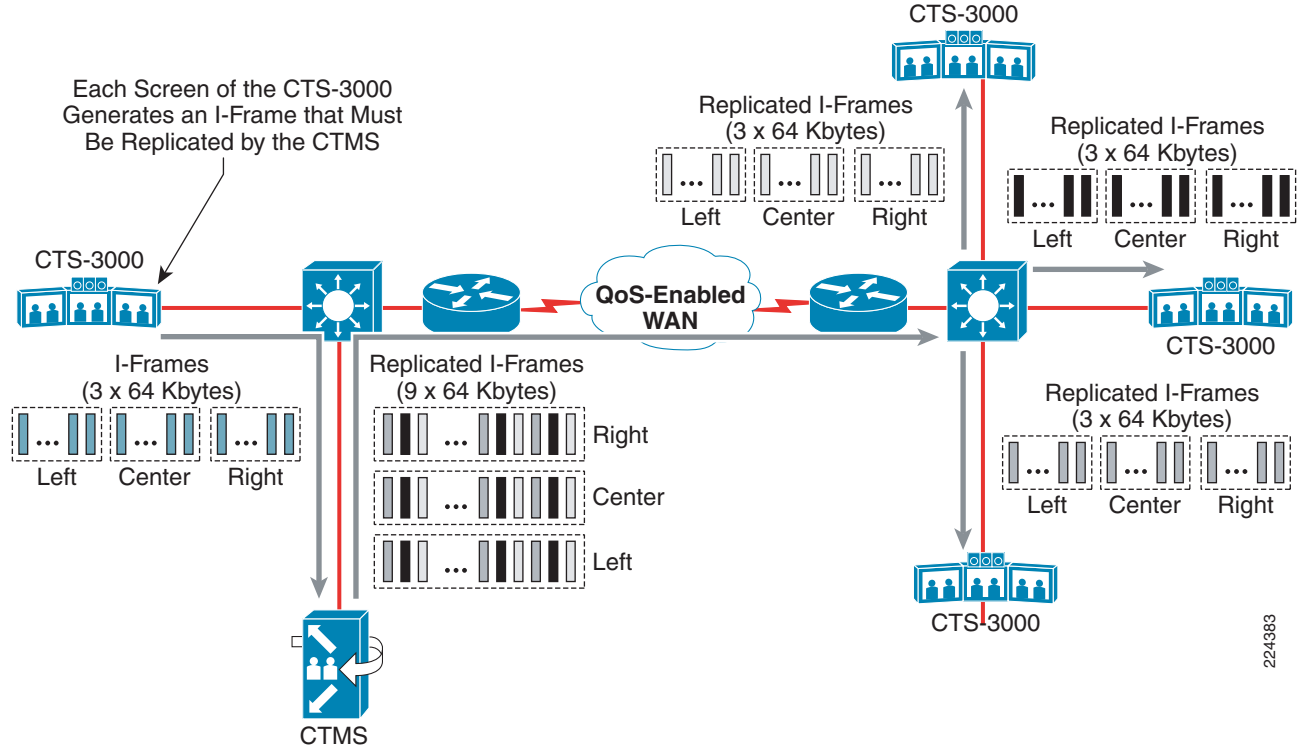
Type of CTS Endpoints

The type of CTS endpoints and whether they are configured for room switching or speaker switching within the multipoint call also determines the size of the bursts generated due to I-frame replication by the CTMS.

Bursts as a Result of Room Switching

[Figure 11-22](#) shows the same four-site multipoint call as shown in [Figure 11-20](#). However, this time the CTS endpoints are CTS-3000 units instead of CTS-1000 units. Further, the multipoint call is configured for room switching.

Figure 11-22 Burst Size Due to CTS-3000 Room Switching

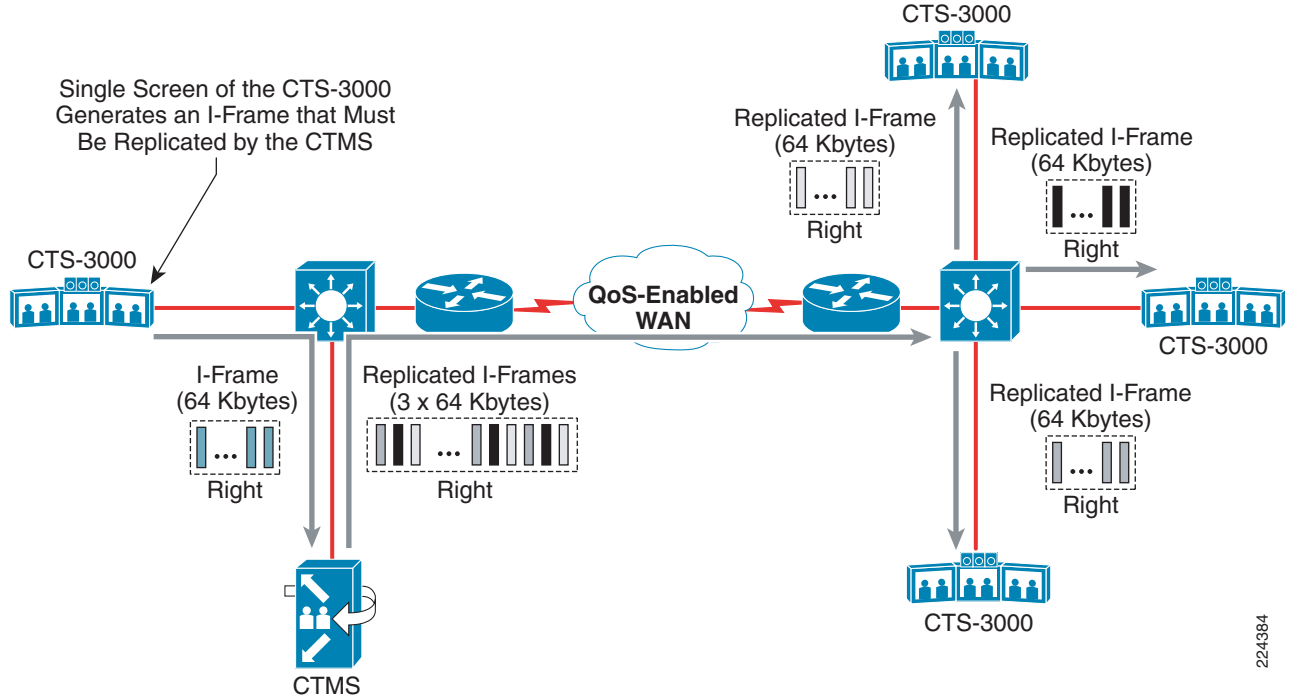


With room switching enabled, whenever any one of the participants at a CTS-3000 endpoint talks and becomes the active speaker, video from all three screens from the CTS-3000 is transmitted to every other CTS-3000. This means that three I-frames are generated, one for each screen position (left, center, and right). These are each replicated by the CTMS to each of the other rooms in the multipoint call. In the example above, the size of the I-frame bursts sent across the WAN are now approximately 3 sites x 64 Kbytes x 3 screens = 576 Kbytes.

CTS-3000s do slightly stagger the I-frame generation since the codecs each handle video independently. Therefore, although each I-frame is generated within a single 33 ms interval, the total duration of the sum of the three I-frames may extend across more than 33 ms. However, any policer and/or shaper parameters configured on routers or within the WAN which have a time constant (otherwise known as the refresh interval) greater than one frame interval ($T_c > 33$ ms), may need to be configured to handle the entire amount of I-frame bursts sent by the CTMS due to room switching.

Bursts as a Result of Speaker Switching

Figure 11-23 shows the same four-site multipoint call again. However, this time the CTS-3000 endpoints are configured for speaker switching.

Figure 11-23 Burst Size Due to CTS-3000 Speaker Switching

With speaker switching enabled, whenever any one of the participants at a CTS-3000 site talks and becomes the active speaker, only the video from the one segment is transmitted to every other participant. This behavior is similar to CTS-1000 endpoints as shown in [Figure 11-20](#). In the example above, only the right screen position sends an I-frame. The size of the I-frame bursts replicated by the CTMS and sent across the WAN is again 3 sites x 64 Kbytes x 1 screen = 192 Kbytes. Combinations of CTS-1000s and CTS-3000s or CTS-3200s in a single call configured for speaker switching behave in a similar manner. Therefore room switching produces more burstiness across the network, although the I-frames may be slightly staggered by the CTS-3000 and CTS-3200 endpoints.

Calculating Burst Sizes Due to I-Frame Replication

Calls with CTS-1000s Only

An estimation of the maximum size of the burst replicated by the CTMS due to I-frames generated during an event such as a normal speaker transition, periodic synchronization, or synchronization due to loss, can be estimated as follows for calls with CTS-1000s only:

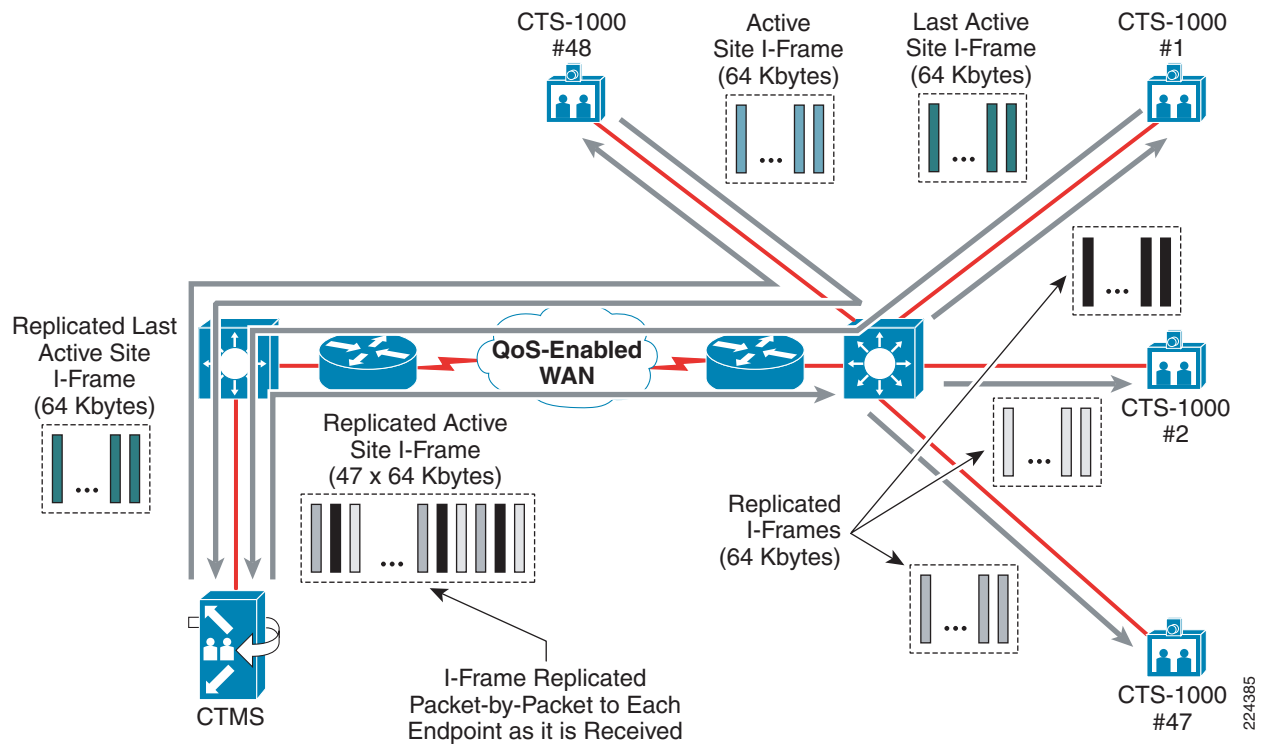
$$\text{CTMS replicated burst size} = (N - 1) * 64 \text{ Kbytes}$$

Where N is the number of CTS-1000s in the multipoint call.

Note that during these events, the number of I-frames replicated by the CTMS is one less than the number of CTS-1000 endpoints within the multipoint call.

However, a worst case scenario for a burst due to I-frame replication by the CTMS occurs when 48 CTS-1000s are in a single multipoint call, the call is configured with the Video Announce feature, and the last CTS-1000 endpoint joins the call. An example of this is shown in [Figure 11-24](#).

Figure 11-24 Maximum Burst Size Generated by the CTMS Due to I-Frame Replication



In such a configuration the size of the I-frame burst replicated by the CTMS is:

$$\text{CTMS replicated burst size} = N * 64 \text{ Kbytes}$$

Where N is the number of CTS-1000s in the multipoint call.

This can be calculated from the example above to be:

$$48 \text{ CTS-1000s} * 64 \text{ Kbytes} = 3.072 \text{ Mbytes}$$

When the CTS-1000 #48 joins the call and becomes the active site, it generates an I-frame which is replicated by the CTMS to the 47 other sites as shown in Figure 11-24. The last active site, CTS-1000 #1 in the example above, also generates a new I-frame which is replicated by the CTMS and sent to CTS-1000 #48 for it to display on its screen. Therefore, in this situation, the total number of I-frames replicated by the CTMS is equal to the number of endpoints in the call. However, it should be noted that the I-frame generated by the last active speaker is not time synchronized with the rest of the I-frames generated by the new CTS-1000 joining the call. In other words, this I-frame may not occur within the same 33 ms frame window as the other 47 I-frames. However, any policer and/or shaper parameters configured on routers or within the WAN, which have a refresh interval greater than one frame interval ($T_c > 33 \text{ ms}$), may need to be configured to handle the entire amount of I-frame bursts sent by the CTMS in this scenario.

Calls with CTS-3000s or CTS-3200s Only

An estimation of the maximum size of the burst replicated by the CTMS due to I-frames generated during an event such as a normal speaker transition, periodic synchronization, or synchronization due to loss can be estimated as follows for calls with CTS-3000s or CTS-3200s only:

$$\text{CTMS replicated burst size} = 3 \text{ screens} * (M - 1) * 64 \text{ Kbytes}$$

Where M is the number of CTS-3000s or CTS-3200s in the multipoint call.

Keep in mind that since CTS-3000 and CTS-3200s stagger their I-frame generation slightly, this burst may occur over more than one 33 ms interval.

As in the previous section, a worst case scenario burst due to I-frame replication by the CTMS occurs when 16 CTS-3000s or CTS-3200s are in a single multipoint call, the call is configured for room switching, the Video Announce feature is enabled, and the last CTS-3000 or CTS-3200 endpoint joins the call. In this case, the size of the replicated I-frame burst by the CTMS is given by:

$$3 \text{ screens} * M * 64 \text{ Kbytes}$$

Where M is the number of CTS-3000s or CTS-3200s in the multipoint call.

Therefore, in this situation, the total number of I-frames replicated by the CTMS is equal to three times the number of endpoints in the call (because the CTS-3000s and CTS-3200s have three screens). For the specific case of 16 CTS-3000s this results in the following:

$$3 * (16 \text{ CTS-3000s}) * 64 \text{ Kbytes} = 3.072 \text{ Mbytes}$$

However, it should be noted that the I-frame generated by the last active site is again not time synchronized with the rest of the I-frames generated by the new CTS-3000 joining the call.

Calls with Mixed CTS-1000s and CTS-3000s or CTS-3200s

The maximum I-frame bursts which result from speaker transitions, periodic synchronization, or synchronization due to loss during mixed CTS-1000 and CTS-3000 or CTS-3200 calls, occurs when the CTS-3000s or CTS-3200s are configured for room switching and a CTS-3000 or CTS-3200 becomes the active site. An estimation of the maximum size of the burst replicated by the CTMS due to I-frames generated during such an event can be estimated as:

$$\text{CTMS replicated burst size} = (3 * (M-1) + N) * 64 \text{ Kbytes}$$

Where N is the number of CTS-1000s in the multipoint call and M is the number of CTS-3000s or CTS-3200s in the multipoint call.

Again, keep in mind that since CTS-3000s and CTS-3200s stagger their I-frame generation slightly, this burst may occur over more than one 33 ms interval.

As with the previous two sections, a worst case burst due to I-frame replication by the CTMS occurs when the call is configured for room switching, the Video Announce feature is enabled, and the last CTS-3000 or CTS-3200 endpoint joins the call. In this case, the size of the replicated I-frame burst by the CTMS will be given by:

$$((3 * M) + N) * 64 \text{ Kbytes}$$

Where M is the number of CTS-3000s or CTS-3200s and N is the number of CTS-1000s in the multipoint call.

For example, if a single multipoint call had six CTS-3000s and 30 CTS-1000s, the worst case burst due to I-frame replication would occur if a CTS-3000 was the last to join the call and would be as follows:

$$(3 * (6 \text{ CTS-3000s}) + 30 \text{ CTS-1000s}) * 64 \text{ Kbytes} = 3.072 \text{ Mbytes}$$

However, it should be noted that the I-frame generated by the last active speaker is again not time synchronized with the rest of the I-frames generated by the new CTS-3000 joining the call. Also the size of the burst is different if the last device to join the call is a CTS-1000 instead of a CTS-3000.

Other Considerations

I-frame bursts due to speaker transitions and periodic synchronization of the video are normal re-occurring events within a multipoint call. The network should be designed to handle bursts generated by these events. I-frame bursts generated when the last CTS-endpoint joins a call is typically a one-time

event at the start of the meeting. Therefore, the design engineer has some discretion regarding whether to design the network to support the entire burst size generated by such an event, considering that I-frame generated by the last active speaker is not time-synchronized with the other I-frames. Also, the Video Announce feature can be disabled to minimize bursts at the beginning of calls.

Next, the design engineer should keep in mind that every network device may not see the entire burst replicated by the CTMS, depending upon the location of the endpoints on the network. The design engineer should determine the maximum number of I-frames which may traverse a particular network device in order to understand the shaper and policer parameters to configure.

Finally, the design engineer should keep in mind that the size of the I-frames replicated by the CTMS is highly variable. The value of 64 Kbytes represents a maximum size empirically observed by testing in the ESE lab. Actual sizes of I-frames generated by CTS endpoints may be smaller than 64 Kbytes. However, it is advisable to design the network to accommodate the maximum observed I-frame sizes. Failure to accommodate this could result in a I-frame “storm” which ultimately may degrade the video quality or cause the multipoint call to fail.

Bursts due to the Auxiliary Video Input

The size of bursts generated from the auxiliary video input is highly dependent upon the content being displayed. The current low-speed auxiliary video input is limited to approximately 500 Kbps without network overhead (approximately 577 Kbps with network overhead). The frame rate of the auxiliary video input is 5 frames per second.

Continuous motion inputs, such as a video clip or continuous animation on a PowerPoint slide being displayed through the auxiliary video input, tend to produce a relatively high bit rate close to 500 Kbps, with fairly uniform bursts of 3-4 Kbytes every 200 ms. On the other hand, PowerPoint slide transitions tend to produce an overall lower bit rate of around 200 Kbps, but much larger bursts. Bursts as high as 44 Kbytes within a single 200 ms frame interval have been observed. Packet sizes average around 1,100 bytes.

However, due to the upcoming release of high-speed auxiliary video which relies on a separate codec and operates at 4 Mbps with 30 frames per second as well, a recommended approach is to account for the auxiliary video burst in a similar fashion as another I-frame being generated by another camera video input, since it represents a worse case than the existing low-speed auxiliary video input. Therefore a value of 64 Kbytes is utilized for all calculations within this document.

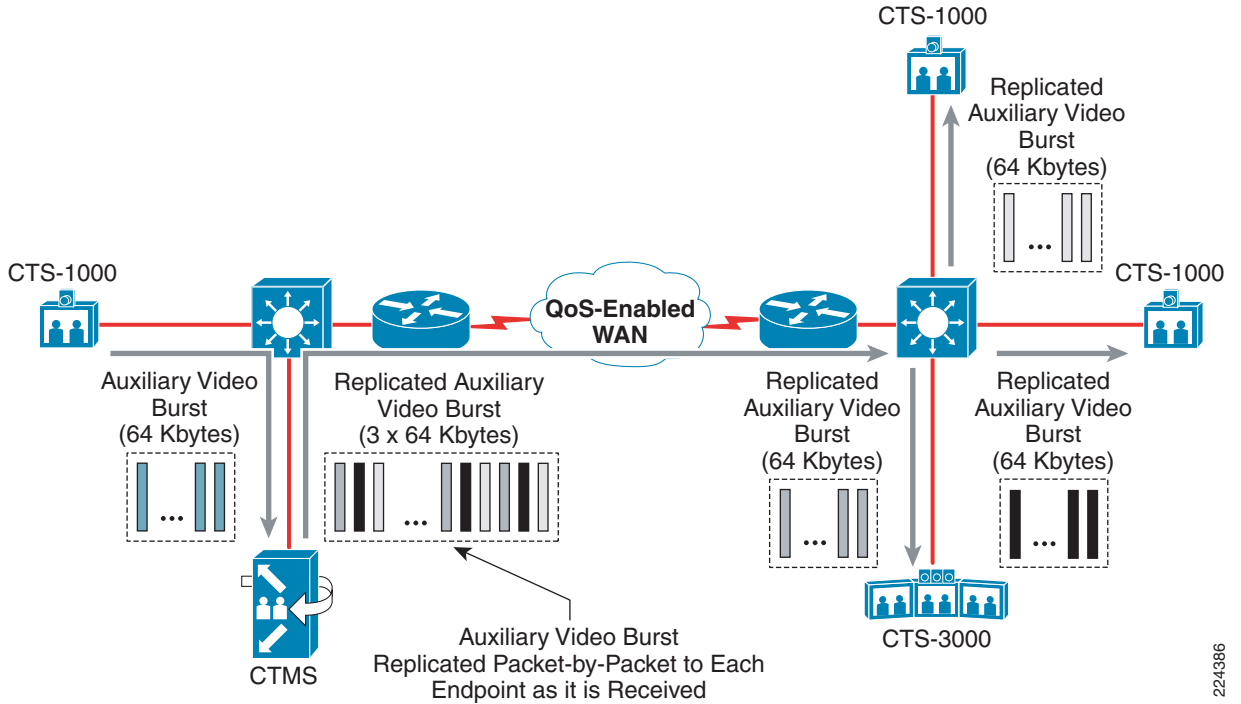


Note

ESE has currently not tested high-speed auxiliary video input in order to confirm that bursts do not exceed 64 Kbytes at this time. Some caution is needed until this can be verified.

As with I-frames from speaker transitions auxiliary video bursts are also replicated on a packet-by-packet basis by the CTMS to every endpoint in the multipoint call, as shown in [Figure 11-25](#).

Figure 11-25 Burst Size Due to Auxiliary Video Burst Replication by the CTMS



It can be seen that as the number of CTS endpoints in the call increases, the size of the replicated auxiliary video burst also increases. However, the type of CTS unit (CTS-1000, CTS-3000, or CTS-3200) does not matter, because each CTS endpoint only has one auxiliary video output.



Note

Simultaneous use of the low-speed and high-speed auxiliary video input in a single multipoint call is not supported.

Burst Estimation Due to Auxiliary Video Replication

An estimation of the maximum size of the burst replicated by the CTMS due to auxiliary video replication can be calculated as:

$$\text{CTMS replicated burst size} = (M + N - 1) * 64 \text{ Kbytes}$$

Where M is the number of CTS-3000 or CTS-3200s and N is the number of CTS-1000s in the multipoint call.

Since one CTS endpoint serves as the presenter, auxiliary video is replicated to one less than the number of CTS endpoints in the multipoint call. In the example above, the size of the auxiliary video burst sent by the CTMS is approximately 3 x 64 Kbytes = 192 Kbytes. This is sent within a single auxiliary video frame interval. For low speed auxiliary video the frame interval is 200 ms, which is a different frame interval than the video from the CTS endpoint cameras. For high speed auxiliary video the frame interval is 33 ms, which is the same frame interval as the video from the CTS endpoint cameras.

**Note**

ESE testing has shown that low-speed auxiliary video content is replicated to every CTS endpoint in the multipoint call, regardless of whether the CTS endpoint is configured to support a projector output within CUCM or whether a projector is actually connected to the output of the CTS endpoint. Note that the behavior of high-speed auxiliary video input has not been tested by ESE currently.

The maximum auxiliary video burst from the CTMS occurs with 48 CTS-1000s in a single multipoint call. An estimation of the maximum size of the burst generated during such an event can be estimated as:

$$\text{CTMS replicated burst size} = (48 - 1) * 64 \text{ Kbytes} = 3.008 \text{ Mbytes}$$

Other Considerations

Auxiliary video bursts due to PowerPoint slide transitions are normal re-occurring events within a multipoint call which utilizes the Auto-Collaborate feature. The network should be designed to handle bursts generated by these events.

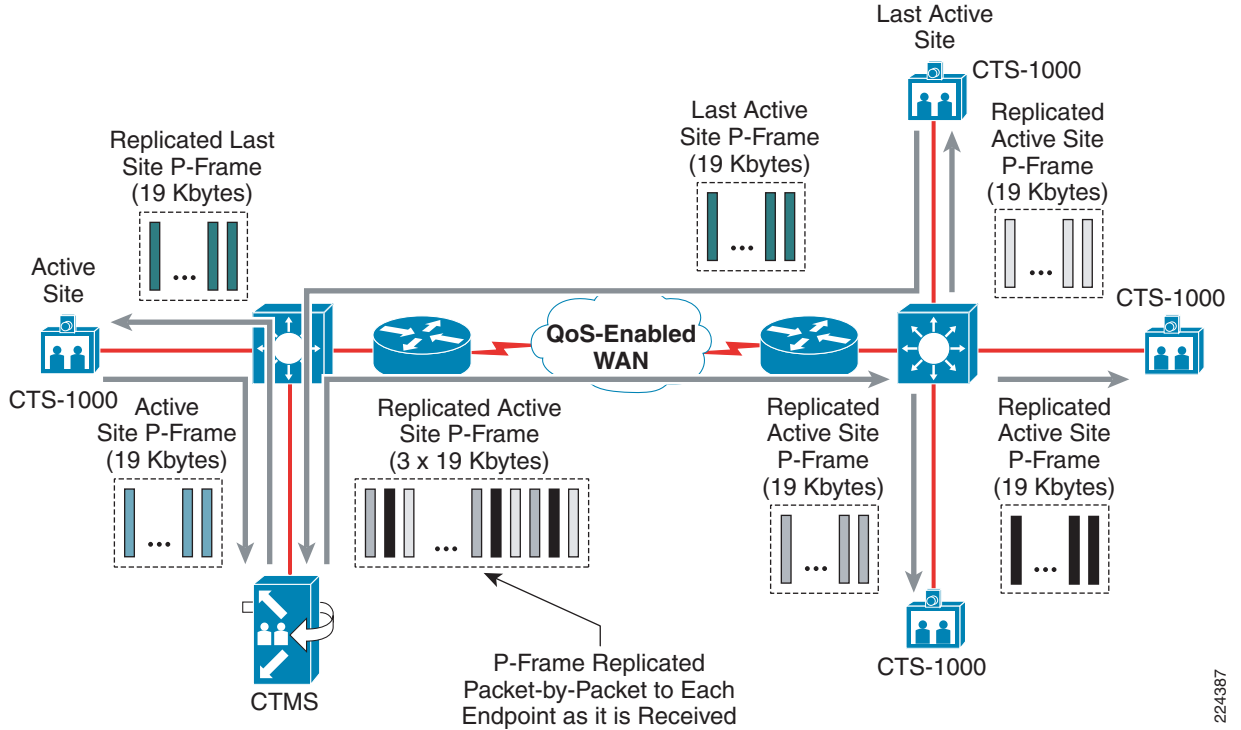
As with I-frame bursts, the design engineer should keep in mind that every network device may not see the entire burst replicated by the CTMS, depending upon the location of the endpoints on the network. The design engineer should determine the maximum number of auxiliary video bursts which may traverse a particular network device in order to understand the shaper and policer parameters to configure.

Finally, the design engineer should keep in mind that the size of the auxiliary video bursts replicated by the CTMS is highly variable. A value of 44 Kbytes over a 200 ms frame interval has been confirmed from ESE testing of low-speed auxiliary video input using PowerPoint slides. However, with the upcoming release of high-speed auxiliary video, a recommended value of 64 Kbytes (matching existing codec I-frame bursts) is recommended for design purposes. Note that high-speed auxiliary video operates over a 33 ms frame interval as well.

Normal P-Frame Video

When the codecs of a CTS endpoint are not sending I-frames to re-synchronize the video, they send P-frames every 33 ms. The size of the P-frames is variable and depends on how well motion estimation and compensation was able to compress the frame. The video from each camera connected to the CTS codec is constrained to approximately 4 Mbps without network overhead or approximately 4.616 Mbps with network overhead. Without any I-frames, this translates to a maximum of approximately 19 Kbytes of video sent every 33 ms, although typically video in a TelePresence meeting runs below this value. The P-Frame video is also replicated by the CTMS on a packet-by-packet basis as shown in [Figure 11-26](#).

Figure 11-26 Burst Size Due to P-Frame Replication by the CTMS



As can be seen from Figure 11-26, the active site generates P-frames which are sent to the CTMS and replicated to the other CTS endpoints. However, since the active site also has to display video, the last active site generates P-frames which are replicated by the CTMS and sent to the active site. In the example above, the size of the P-frame bursts sent by the CTMS is approximately $4 \times 19 \text{ Kbytes} = 76 \text{ Kbytes}$ every frame interval. Note however that this fourth P-frame is not time synchronized (does not happen exactly at the same time) with the rest of the P-frames replicated by the CTMS.

Location of the CTS Endpoints

As with I-frame bursts, the location of the CTS endpoints within the multipoint call influences the size of the replicated P-frame burst across a given network device or WAN circuit. In the example above, four P-frames are replicated by the CTMS (three from the active site and one from the last active site) and seen by the switch directly connected to the CTMS. However only three P-frames or $3 \times 19 \text{ Kbytes} = 57 \text{ Kbytes}$ are sent across the WAN to the remote CTS endpoints every 33 ms.

Burst Estimation Due to P-Frame Replication

The type of CTS endpoints within the multipoint call affects how the P-frame bursts are generated, but not the size of the overall bursts. For meetings which include CTS-3000s or CTS-3200s in speaker switching mode, or combinations of CTS-3000s, CTS-3200s, and CTS-1000s, there may be multiple last active segments since each screen of each CTS-3000 or CTS-3200 needs to display video. For meetings which include only CTS-3000s or CTS-3200s in room switching mode, the active site sends three P-frames every 33 ms, one from each codec. Likewise the last active site sends three P-frames every 33 ms which is replicated by the CTMS and sent to the active speaker.

However, in all cases, the number of P-Frames replicated by the CTMS during a frame interval, regardless of the type of CTS endpoint and whether it is configured for room or speaker switching, is given by the following equation:

Number of P-Frames Replicated by the CTMS = Number of Video Segments in the Multipoint Call

Therefore, as the number of CTS endpoints in the call increases, the size of the replicated P-frame burst also increases.

An estimation of the worst case scenario of a P-frame burst generated by the CTMS per frame interval, when 48 video segments are in a single multipoint , is:

CTMS replicated burst size = $48 * 19 \text{ Kbytes} = 912 \text{ Kbytes}$

Other Considerations

P-frame bursts are normal re-occurring events within a multipoint call. The network must be designed to handle bursts generated by these events. However, the design engineer should keep in mind that every network device may not see the entire burst replicated by the CTMS, depending upon the location of the endpoints on the network. The design engineer should determine the maximum number of P-frames which may traverse a particular network device in order to understand any shaper and policer parameters to configure.

Finally, the design engineer should keep in mind that the size of the P-frames replicated by the CTMS is highly variable as well. The value of 19 Kbytes represents an average size based upon the maximum video rate of 4 Mbps (4.616 Mbps with network overhead). Normal video is typically below this rate, suggesting smaller P-frame bursts. However, it is advisable to design the network to accommodate the P-frame bursts based on maximum video rate, since this represents a worst case. Failure to accommodate P-frame bursts causes the video quality to degrade and the multipoint call to fail.



CHAPTER 12

Cisco TelePresence Multipoint Solution Circuit and Platform Recommendations

Recommendations for Multipoint over WAN Circuits

This section discusses the recommendations for support of multipoint and multiple point-to-point TelePresence over private WAN circuits. Recommendations are made with respect to supporting TelePresence in a dedicated circuit configuration and a converged circuit configuration. Dedicated circuit configurations are those in which TelePresence is the only traffic which exists on the private WAN circuit. In other words, the customer has deployed an overlay network designed exclusively for the support of TelePresence. Converged circuit configurations are those in which TelePresence is integrated with data, voice, and other video applications over a private WAN infrastructure.



Note

Recommendations for Multipoint over WAN circuits which involve a handoff to a service provider network, such as MPLS or Metro-Ethernet, and which involve shaping to limit the data rate to the service provider, are discussed in [Multipoint TelePresence over MPLS Circuits with Ethernet Handoff](#).

Cisco Router and Switch Platforms Tested

The following router and switch platforms were tested for the recommendations which follow:

WAN Catalyst 6500 switches:

- Catalyst 6506-E
 - IOS Version 12.2(18)SXF12
 - WS-SUP32-10GE-3B
 - WS-F6K-PFC3B
 - WS-F6K-MSFC2A
- 7600 SIP-400
 - SPA-2XOC3-POS
 - SPA-1XOC12-POS
 - SPA-1XOC48POS/RPR
- 7600-SIP-200
 - SPA-4XT3/E3

WAN Cisco 7200 routers:

- Cisco 7206VXR (NPE-G2) processor (rev A) with 917504K/65536K bytes of memory.
- IOS Version 12.4(4)XD7
- PA-POS-OC3MM
- PA-T3+

TelePresence over Dedicated WAN Circuits

The following sections discuss the recommendations for support of multipoint and multiple point-to-point TelePresence over dedicated private WAN circuits.

Dedicated T3, E3, and OC-3 POS Circuits

Table 12-1 summarizes the recommendations for supporting multipoint TelePresence over T3, E3, and OC-3 WAN circuits.

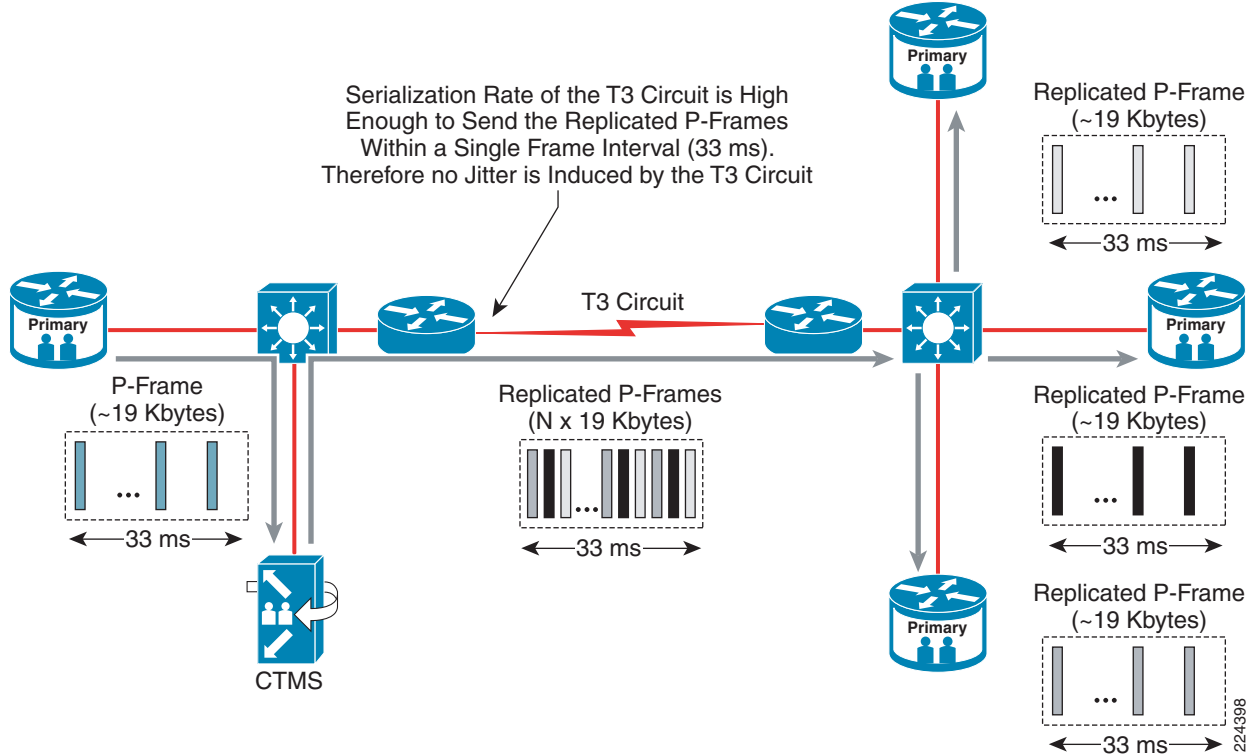
Table 12-1 Multipoint over Dedicated T3, E3, and OC-3 Circuit Results

Circuit	Bandwidth	Maximum Recommended CTS Table Segments	Recommended Starting Point for Tuning Output Hold Queue	TxRing Settings
T3	44.2 Mbps	7	840 (See discussion)	Default (See discussion)
E3	34.01 Mbps	5	600 (See discussion)	Default (See discussion)
OC-3 POS	155 Mbps	22 Estimated (See discussion)	2,640 (See discussion)	Default (See discussion)

Multipoint Recommendations

When configured in a single multipoint call, the maximum number of CTS table segments for lower speed WAN circuits, such as T3, E3, and OC-3 POS, is bounded not by the bandwidth of the circuits, but by the jitter induced by serialization delay across the circuits. The recommended maximum number of CTS table segments in a single multipoint call over such circuits is based upon the number of I-Frames and/or auxiliary video (PowerPoint slide) bursts simultaneously replicated by the CTMS and sent across the circuit. Figure 12-1 and Figure 12-2 help explain this using a T3 circuit as an example.

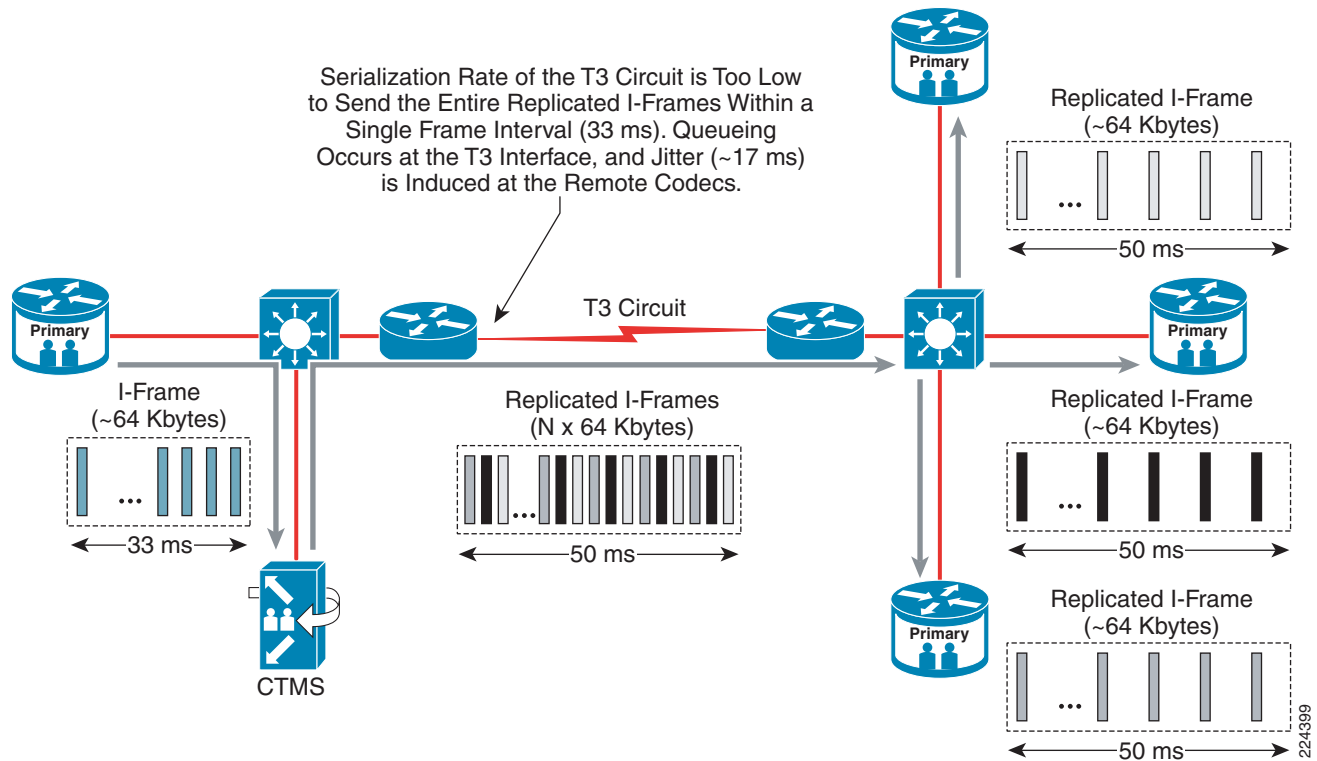
Figure 12-1 Example P-Frame Replication Across a T3 Circuit



As shown in [Figure 12-1](#), the transmitting codec sends a video frame from its camera every 33 ms frame interval. In this case, the video frame is a highly-compressed P-frame. The P-frame is replicated by the CTMS to the other three remote codecs across the T3 circuit. The replicated information is three times the size of the original P-frame. However, it must be received by each remote codec within the same 33 ms frame interval. In the example in [Figure 12-1](#), the serialization rate of the T3 circuit is high enough such that the entire replicated P-frames can still be sent across the circuit within 33 ms. Therefore, each remote codec receives its P-frame within the frame interval and no jitter is induced. This assumes that the latency due to physical distance of the circuit is constant.

[Figure 12-2](#) shows the same configuration with a I-frame replicated across the T3 circuit.

Figure 12-2 Example I-Frame Replication Across a T3 Circuit



As shown in Figure 12-2, the transmitting codec now sends a I-frame from its camera during the 33 ms frame interval. I-frames are about 3 to 4 times larger than P-frames, so there is much more information to be sent. The I-frame is again replicated by the CTMS to the other three remote codecs across the T3 circuit. The replicated information is three times the size of the original I-frame. It must still be received by each remote codec within the same 33 ms frame interval. In the example in Figure 12-2, the serialization rate of the T3 circuit is not fast enough to send the entire replicated I-frame within 33 ms. Therefore, queuing occurs at the T3 interface of the sending router. In the example above, the entire replicated I-frames are sent out in 50 ms. Since the CTMS replicates on a packet-by-packet basis, each remote CTS codec receives the full I-frame over a 50 ms interval. Therefore, approximately 50 ms - 33 ms = 17 ms of jitter is induced in each of the remote codecs when the I-frame is sent.

The tolerance of the CTS endpoints to jitter is based upon continuous jitter over a 50 second time period. Since jitter induced by an I-frame or auxiliary PowerPoint burst is a transient event, there is no concern that the call will terminate due to jitter exceeding the target of 40 ms for bad video quality. Also, the transient jitter can exceed the target of 20 ms for good video quality without causing any visible artifacts on the displays. However, the de-jitter buffers of the codecs are somewhat dynamic, with a target depth somewhere between 60-80 ms. If the transient jitter induced by I-frames and auxiliary bursts (due to PowerPoint slide transitions) exceeds the current depth of the de-jitter buffer, the codecs will mark packets as being "late." Late packets are not used in the decoding process of the codecs and cause visible artifacts on the displays. This degrades the video quality and the overall TelePresence experience.

Based upon ESE testing with CTS-1000s in a single multipoint call with low-speed (5 frames per second) auxiliary video input, the maximum number of endpoints which were observed to be supported on T3 and E3 circuits before late packets were reported on the CTS codecs is listed in Table 12-1. Note that these recommendations are based on testing across a single circuit.

**Note**

The number of CTS segments which can be supported with high-speed auxiliary video input may be lower due to the higher speed of the auxiliary video input and potentially larger bursts. ESE has not tested high-speed auxiliary video input currently.

The maximum number of segments supported over an OC-3 circuit is estimated based on the results from the T3 and E3 testing, as well as a tested configuration of 15 CTS-1000s and 3 CTS-3000s. Note that the number of replicated I-frames and auxiliary bursts in the tested configuration is somewhat lower than the number that would be generated by 22 CTS-1000s.

Multiple Point-to-Point Recommendations

For circuits which support multiple point-to-point TelePresence calls, the I-frames generated during each call are not time synchronized. In other words, there is no CTMS replication of I-frames to each site simultaneously. As a result of this, less jitter is seen by the codecs and the number of CTS table segments which can be supported across a circuit may be slightly higher. The maximum table segments which may be supported in this configuration should be calculated based on the bandwidth of the circuit and the bandwidth of each CTS endpoint—for example, using 15 Mbps per CTS-3000 and 5.5 Mbps per CTS-1000. Although these values represent maximum rates over a one-second period, it is not recommended to oversubscribe the circuit by provisioning CTS units such that the sum of their one second maximum rates exceeds the bandwidth of the circuit. Additionally, a wise rule-of-thumb for any network is to not exceed 75% utilization of any WAN link on a consistent basis, before considering increasing the bandwidth of the link.

Since the customer may at any time deploy a CTMS and begin supporting multipoint calls across the network, another wise rule-of-thumb is to design the network in such a manner that it is ready for multipoint. Therefore, it is recommended that regardless of whether the customer currently supports multipoint, the network should be designed in such a manner that it can be added without having to adjust network parameters on all circuits throughout the network when multipoint is supported.

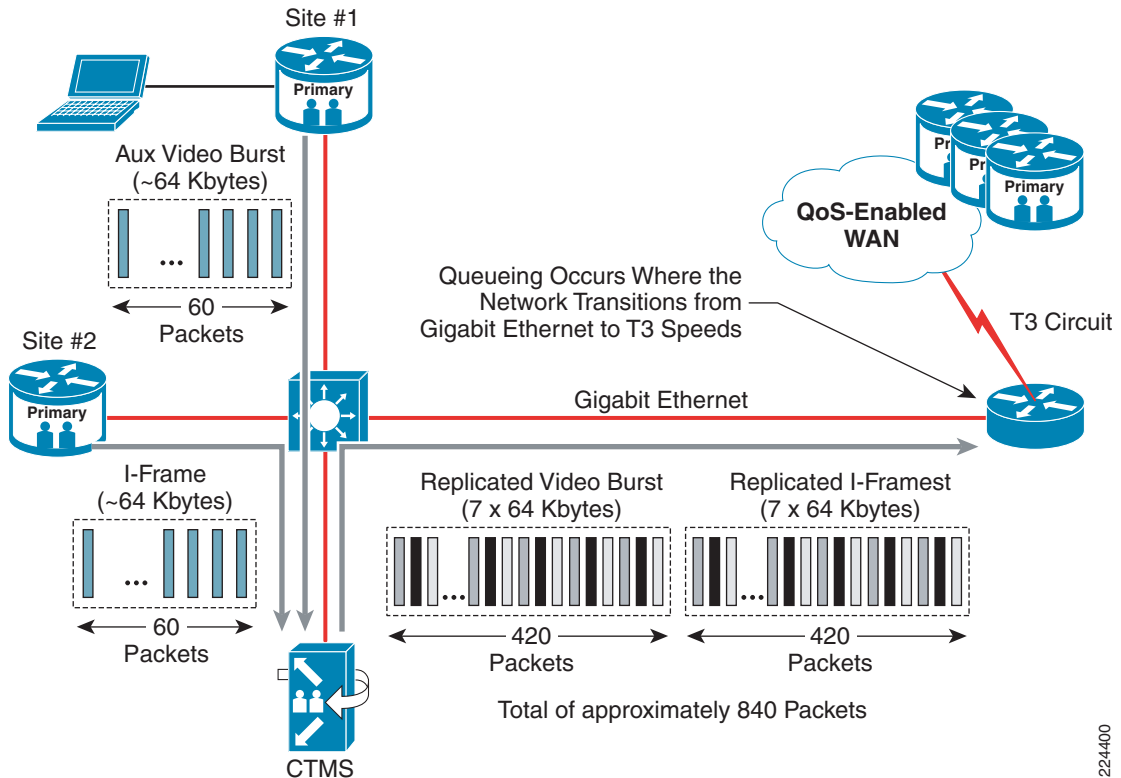
Output Hold Queue Recommendations

In multipoint configurations over dedicated circuits, the output hold queue of routers connected to the circuits may have to be increased to support the bursts due to I-frame replication and auxiliary bursts due to PowerPoint slide transitions.

Buffer size tuning is dependant upon size of the I-frame and auxiliary video (PowerPoint slide) bursts and their timing. Therefore, the output queue sizes listed in [Table 12-1](#) represent starting points for tuning the hold-queue out parameters on router WAN interfaces based upon testing within the ESE lab.

The following example demonstrates one simple method of estimating the output hold queue size based upon information in [Estimating Burst Sizes within Multipoint TelePresence Calls](#) in [Chapter 11](#), “[Cisco Multipoint Technology and Design Details](#).” Assume a multipoint TelePresence call consisting of 9 CTS-1000s. A maximum of seven of the CTS-1000s, totaling seven CTS table segments, are located across a T3 circuit from the CTMS, as shown in [Figure 12-3](#).

Figure 12-3 Output Hold Queue Estimation



224400

In the example in Figure 12-3, a slide show is being presented by Site #1. As the presenter transitions the slide, a question is being asked from Site #2, causing a speaker transition to occur simultaneously. As a result of the two simultaneous events, both an I-frame and an auxiliary video burst are sent simultaneously. Note that audio has been ignored to simplify the example. The I-frame and auxiliary video burst are replicated by the CTMS and sent to the seven CTS 1000s across the T3. (For simplicity of the example, replication of the I-frame and auxiliary video burst to the endpoints on the same side of the T3 circuit is not shown.)

Queuing primarily occurs at the transition from the higher speed Gigabit Ethernet interface to the lower speed T3 interface. In other words, queuing occurs at the output hold queue of the router's T3 interface. The following equations estimate the total burst size in packets:

60 Packets per I-frame * Replication to 7 Sites = 420 Packets

60 Packets per Aux Video Burst * Replication to 7 Sites = 420 Packets

Total Packets = 420 + 420 = 840 Packets

Based upon the example, a good starting point for tuning the output hold queue of the router's T3 interface would be approximately 840 packets. Actual test results for output hold queue sizes required for support of multipoint TelePresence over dedicated circuits have shown the equation above to be slightly conservative. The example above over-simplifies the arrival and transmission of the burst. It assumes that all 840 packets arrive at the same time. In reality, the I-frame is sent over a 33 ms frame interval and the auxiliary video burst is sent over a 200 ms frame interval with low-speed auxiliary video. Also, not shown is a single audio packet which arrives every 20 ms. Likewise, the example assumes that the entire burst arrives before any of it is serialized and sent onto the T3 circuit. In reality, as the information arrives, it is serialized. The higher the speed of the circuit, the more is serialized as it arrives,

and therefore the lower the maximum queue depth. It must also be noted that the size of I-frames and auxiliary video bursts is variable. The values shown above represent estimates for maximum sizes seen during testing.

Despite the over-simplification of the example, the method presented can still be used as a starting point for tuning output hold queues for supporting multipoint TelePresence. More generically, the starting point for tuning the output hold queue can be determined by the equation:

$N * 120 \text{ packets} + M * 240 \text{ packets} = \text{Starting point for tuning the output hold queue size for a dedicated circuit}$

The application of detailed queuing theory to obtain a more accurate assessment of the maximum queue depth is beyond the scope of this document. However, it should be noted that there is no apparent downside, in terms of TelePresence video quality, of simply increasing the output hold-queue on the router to its maximum value. TelePresence simply utilizes as much of the output queue space as needed. Bounding the maximum number of CTS table segments supported across the circuit and then allowing as much queuing space as needed is preferable to discarding TelePresence traffic.

TxRing Tuning Recommendations

For multipoint TelePresence over a dedicated circuit, it is recommended to leave the TxRing configuration of the WAN interface at its default value. Testing over a T3 circuit has indicated that decreasing the value to 10 packets—as was the recommendation for point-to-point TelePresence over a shared T3 circuit—resulted in output drops on the T3 interface. The length of the TxRing determines when queuing is enabled and disabled due to congestion on IOS routers. In a dedicated circuit configuration, TelePresence is the only traffic. Therefore, queueing is somewhat irrelevant and the TxRing can be left at its default value. Note that within many POS interfaces on IOS routers, the TxRing cannot be tuned anyway.

Multipoint Over Dedicated OC-12 and OC-48 POS Circuits

Testing the maximum number of CTS table segments which a dedicated OC-12 or OC-48 POS circuit is capable of handling is beyond the capability of the ESE lab. However both OC-12 and OC-48 POS have been tested with 15 CTS-1000s and 3 CTS-3000s in a single multipoint call with low-speed (5 frames per second) auxiliary video input. Additionally, another 8 simulated CTS-3000s were added to each circuit, for a total of 48 CTS table segments. It has been confirmed that dedicated OC-12 and OC-48 POS circuits are capable of handling the equivalent to the amount of TelePresence traffic that a single CTMS can currently generate across these circuits.

No tuning of output hold queues was necessary for the tested configuration on OC-12 or OC-48 circuits. Further, as mentioned previously, the TxRing cannot be tuned for the tested OC-12 and OC-48 interfaces.

Recommendations for Multipoint Over Converged WAN Circuits

The following sections discuss the recommendations for support of multipoint and multiple point-to-point TelePresence over converged WAN circuits. Converged circuit configurations are those in which TelePresence is integrated with data, voice, and other video applications over a private WAN infrastructure.

When configuring multipoint over a WAN edge design, the network administrator should follow the basic design guidance outlined in [Chapter 6, “Branch QoS Design for TelePresence.”](#) The following sections extend those discussions to multipoint TelePresence designs.

Multipoint TelePresence WAN Edge LLQ Policy

When configuring multipoint TelePresence over a converged private WAN infrastructure, it is often advantageous to place TelePresence in an LLQ. If multipoint TelePresence is to be assigned to an LLQ, then sufficient bandwidth must be allocated to the LLQ to support multiple CTS endpoints. The LLQ bandwidth is typically limited by either implicit policing or a configured explicit policer. In either case, the policed rate (CIR) must be sufficient for the number of CTS endpoints supported across the converged WAN circuit. This is regardless of whether the CTS endpoints are in a single multipoint call, multiple multipoint calls, or multiple point-to-point calls.

Bandwidth Provisioning

The amount of LLQ bandwidth provisioned via the policed rate can be calculated based upon the number of CTS endpoints to be supported across the WAN circuit and the resolution/quality configuration for the endpoints, as shown in [Table 4-1](#) in [Chapter 4](#), “Quality of Service Design for TelePresence.”

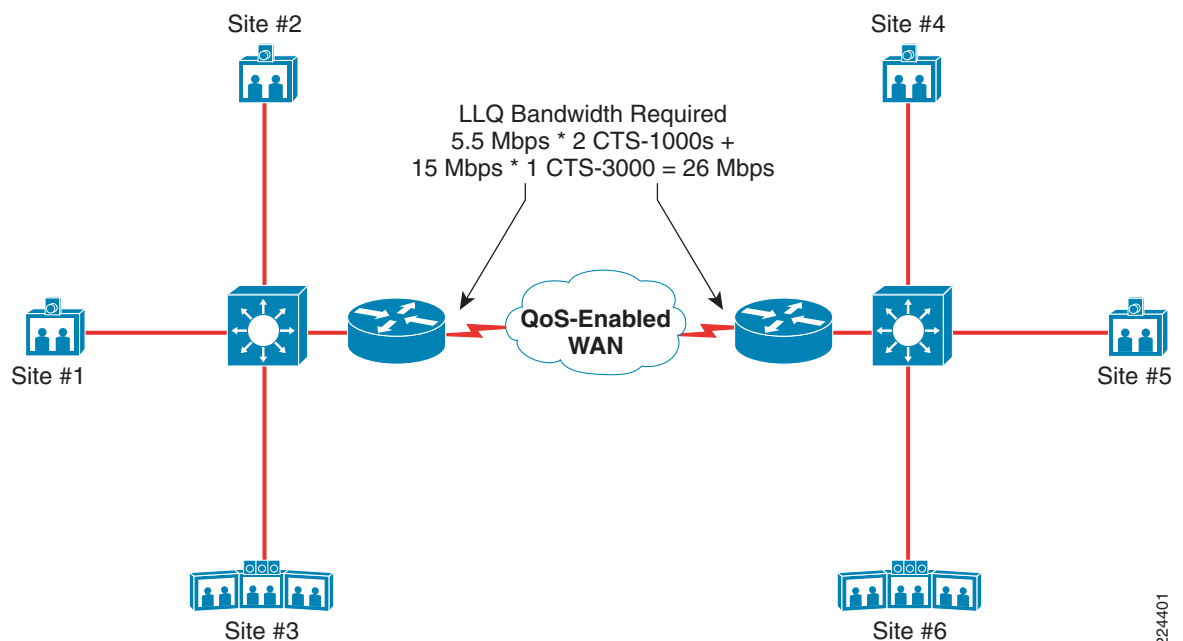
For example, with 1080p best quality the amount of provisioned LLQ bandwidth can be determined using the equation:

$$N * 5.5 \text{ Mbps} + M * 15 \text{ Mbps} = \text{Policed Rate (CIR)}$$

Where N is the number of CTS-1000s and M is the number of CTS-3000s to be supported across the circuit.

[Figure 12-4](#) shows this with an example TelePresence deployment.

Figure 12-4 Bandwidth Requirements for Multipoint and Multiple Point-to-Point TelePresence



In [Figure 12-4](#) there are a total of six CTS endpoints in the network—four CTS-1000s and two CTS-3000s. However each side of the WAN has only three CTS endpoints. The amount of bandwidth required to be provisioned in the LLQ policed rate across the WAN, in order to support having all TelePresence units in calls at the same time, is:

$$2 \text{ CTS-1000s} * 5.5 \text{ Mbps per} + 1 \text{ CTS-3000} * 15 \text{ Mbps} = 26 \text{ Mbps}$$

224401

This holds regardless of whether the CTS units are in a single multipoint call or in multiple multipoint calls.



Note

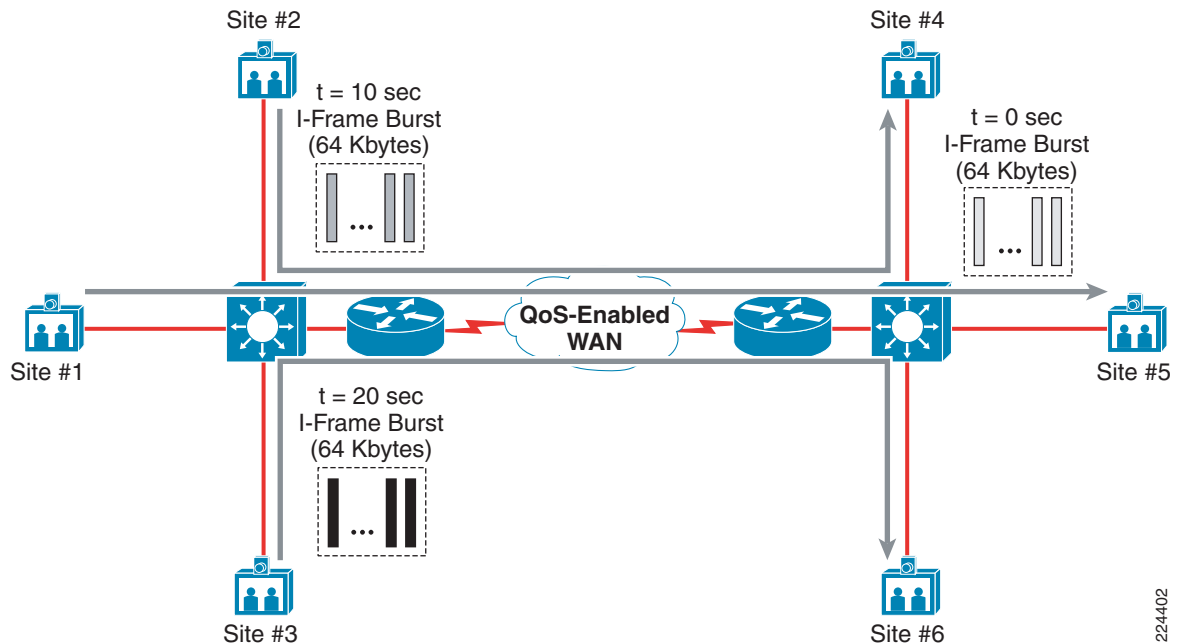
These recommendations hold for low-speed (5 frames/second) auxiliary video input only. Additional bandwidth must be provisioned to support high-speed (30 frames/second) auxiliary video input.

It should be noted that with TelePresence in an LLQ configuration, it is recommended that the total amount of traffic allocated for the LLQ (Voice and TelePresence) be below approximately 1/3 of the circuit bandwidth in order to prevent possible decreases in responsiveness of applications which are not placed into the LLQ.

Burst Provisioning

The burst parameter of the LLQ policer must also be scaled to handle bursts from multiple CTS endpoints. This burst size is not directly dependent upon the resolution/quality configuration of the CTS endpoints. However, the required burst size is dependent upon whether CTS endpoints are in individual point-to-point calls, multiple multipoint calls, or one single multipoint call. [Figure 12-5](#) and [Figure 12-6](#) help to explain this.

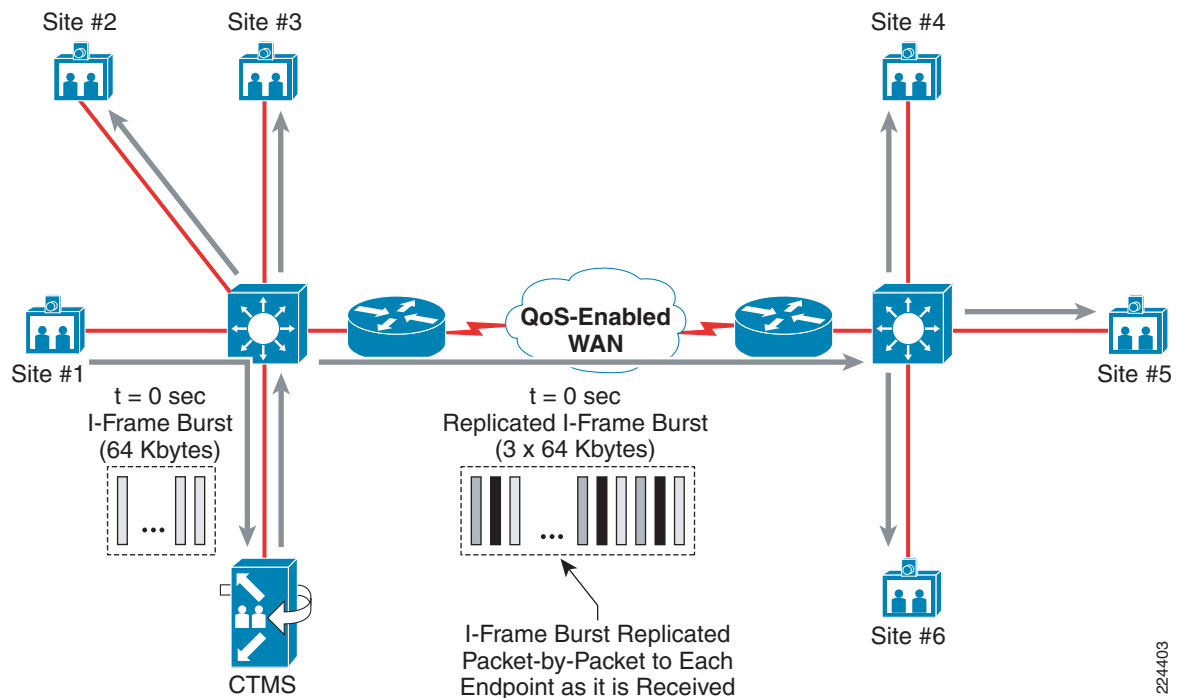
Figure 12-5 Multiple Point-to-Point TelePresence Calls



[Figure 12-5](#) shows an example of six CTS-1000 sites in three point-to-point calls. All three calls traverse the converged WAN. In the example, at time $t = 0$ seconds, Site#1 generates an I-frame burst, perhaps as part of normal video synchronization. Likewise at $t = 10$ seconds, Site #2 generates an I-frame burst, and at $t = 20$ seconds, Site #3 generates an I-frame burst. In the example, the I-frames are not time synchronized. In other words, bursts do not necessarily occur at the same time with separate point-to-point calls. Therefore the burst size of the LLQ policer of the WAN circuit only needs to handle a single I-frame burst in this example.

[Figure 12-6](#) shows the same six sites, in a single multipoint conference call this time.

Figure 12-6 Single Multipoint TelePresence Call



In this example, at time $t = 0$ seconds, Site #1 generates an I-frame burst, perhaps because it became the active site. The I-frame burst is immediately replicated by the CTMS to every other site. Since Sites #4, #5, and #6 are located across the WAN, the LLQ policer burst parameter (B_c) must be configured to handle three times the size of the initial I-frame burst generated by Site #1.

Since a single multipoint call represents a worst case scenario in terms of synchronization of bursts, it is recommended to provision the LLQ burst parameter (B_c) to support this scenario. Allocation of approximately 64 Kbytes of burst space per video camera, as well as another 64 Kbytes of burst space per auxiliary video input on each CTS endpoint, has been shown to be conservatively sufficient. Using this recommendation, each CTS-1000 would require approximately 128 Kbytes of burst space and each CTS-3000 would require approximately 256 Kbytes of burst space.

Using these recommendations, the amount of provisioned LLQ burst for a converged circuit which supports multipoint TelePresence can be calculated based upon the following equation:

$$N * 128 \text{ Kbytes} + M * 256 \text{ Kbytes} = \text{Burst Size } (B_c)$$

Where N is the number of CTS-1000s and M is the number of CTS-3000s to be supported across the circuit.

Going back to the example in Figure 12-3, the burst size required to be provisioned in the LLQ policer in order to support having all TelePresence units in a single multipoint call is:

$$2 \text{ CTS-1000s} * 128 \text{ Kbytes} + 1 \text{ CTS-3000} * 256 \text{ Kbytes} = 512 \text{ Kbytes}$$

This holds regardless of whether the CTS units are in a single multipoint call or in multiple multipoint calls.

Multipoint LLQ Configuration Examples

Example 12-1 and Example 12-2 show both an implicit and explicit policer configuration for support of multipoint TelePresence in an LLQ configuration along with voice.

Example 12-1 Implicit Policer

```

policy-map OC-3 WAN-EDGE
  class VOICE
    priority percent 10! LLQ for Voice (example amount of BW)
  class TELEPRESENCE
    priority 26000 512000! LLQ for TP (priority bandwidth)
  class DATA
  ...

```

Example 12-2 Explicit Policer

```

policy-map OC-3 WAN-EDGE
  class VOICE
    police cir 15500000! Explicit Policer for Voice (default burst)
      conform-action transmit
      exceed-action drop
      violate-action drop
    priority 15500! LLQ for Voice (explicit priority bandwidth)
  class TELEPRESENCE
    police cir 26000000 bc 512000 be 512000! Explicit Policer for TP
      conform-action transmit
      exceed-action drop ! Note: Excess burst (Be) not utilized.
      violate-action drop
    priority 26000! LLQ for TP (explicit priority bandwidth)
  ...

```

It should be noted that the combination of configuring a policed rate of 5.5 Mbps and a burst size 128 Kbytes per CTS-1000 ensures a policer time constant of approximately 186 ms. Likewise, configuring a policed rate of 15 Mbps and a burst size of 256 Kbytes per CTS-3000 ensures a policer time constant of approximately 137 ms. Therefore the recommendations above ensure that the time constant for the policer (given by the equation $T_c = CIR/Bc$) is always between 137 ms and 186 ms. Many service provider networks already utilize time constants of around 200 ms, implying a larger burst value configured for a given policed rate, than recommended above.

Empirical Test Results

Table 12-2 summarizes the test results for supporting multipoint TelePresence over shared T3, E3, and OC-3 WAN circuits, performed within the ESE lab.

Table 12-2 Multipoint over Shared T3, E3, and OC-3 POS Circuit Test Results

Circuit	CTS Unit Configuration	Policed Rate (CIR)	Actual Burst (Bc) Tested	Recommended Burst Size	Percentage of Recommended
E3	2 CTS-1000s	11 Mbps	160 Kbytes	256 Kbytes	63%
T3	3 CTS-1000s	16.5 Mbps	164 Kbytes	384 Kbytes	43%
OC-3	7 CTS-1000s and 1 CTS-3000	53.5 Mbps	680 Kbytes	1,152 Kbytes	59%
OC-3	9 CTS-1000s	49.5 Mbps	720 Kbytes	1,152 Kbytes	63%
OC-12	15 CTS-1000s and 3 CTS-3000s	127.5 Mbps	1,584 Kbytes	2,688 Kbytes	59%

The test results in [Table 12-2](#) show that allocating 128 Kbytes per CTS-1000 and 256 Kbytes per CTS-3000 are fairly conservative recommendations, which should be sufficient for most customer networks. Actual testing has shown that bursts are typically somewhere between 40% to 70% of this value. Several reasons for this wide range of results are:

- The size of I-frames (IDRs) is highly variable, with sizes observed up to 64 Kbytes.
- The size of auxiliary video bursts is also highly variable.
- The probability of part or all of an auxiliary video burst (PowerPoint slide transition) occurring during the exact same time as an I-frame (IDR) is low.

Multipoint TelePresence Branch WAN Edge CBWFQ Policy

If multipoint TelePresence is to be assigned to a CBWFQ, then sufficient bandwidth and queue size for bursts must be allocated to the CBWFQ to support multiple CTS endpoints.

Bandwidth Provisioning

Within a CBWFQ, bandwidth is typically limited via a bandwidth statement within the policy map, as shown in [Example 12-3](#).

Example 12-3 TelePresence CBWFQ Example

```
policy-map T3 WAN-EDGE
  class VOICE
    police cir 1184000! Explicit Policer for Voice (default burst)
      conform-action transmit
      exceed-action drop
      violate-action drop
      priority 1184! LLQ for Voice (explicit priority bandwidth)
  class TELEPRESENCE
    bandwidth percent 37! CBWFQ for TP (example amount of BW)
    queue-limit 120
  ...
```



Note

In [Example 12-3](#), TelePresence traffic can utilize more than 37% of the bandwidth if other CBWFQ classes and/or the voice LLQ class are not utilizing their complete share of the bandwidth of the circuit. This was also the case for TelePresence in the LLQ with an implicit policer.

When allocating CBWFQ bandwidth for multipoint TelePresence, it is recommended to again allocate approximately 5.5 Mbps per CTS-1000 and 15 Mbps per CTS-3000 (assuming low-speed auxiliary video). In the example above which supports three CTS-1000s, bandwidth is calculated as:

$$3 \text{ CTS-1000s} * 5.5 \text{ Mbps per CTS-1000} = 16.5 \text{ Mbps}$$

$$\text{Bandwidth} = 16.5 \text{ Mbps} / 44.21 \text{ Mbps T3 circuit bandwidth} = 37\%$$

Burst Provisioning

When configuring multipoint TelePresence in a CBWFQ, bursting is accommodated by a combination of the queue size configured within the policy map via the queue-limit command, as well as the output hold queue size of the WAN interface configured via the hold-queue out command.

Policy map Queue Size Recommendations

When multipoint TelePresence is placed within a CBWFQ on a converged circuit, the queue size for the TelePresence class within the policy map may need to be increased above the default limit of 64 packets to support the bursts.

There is no exact science as to the size of the queue limit to configure within the policy map for the TelePresence class based upon the number of CTS endpoints supported over the circuit. Tuning of the queue size is dependant upon the sizes of the I-frame and auxiliary video bursts and their timing, as well as the amount of congestion on the WAN circuit which causes queuing to occur in the first place on a converged circuit. Further, TelePresence bursts themselves may be the cause of the temporary congestion which results in queuing.

However, the method shown in [Figure 12-3](#) in [Output Hold Queue Recommendations](#) can be applied to multipoint TelePresence in a CBWFQ in order to provide the network administrator a starting point which to begin tuning the queue size for the TelePresence class.

As shown in [Figure 12-3](#), each I-frame from the camera of a CTS-endpoint can be as large as 64 Kbytes. Given a frame size of approximately 1,100 bytes, this translates to approximately 60 packets per I-frame per camera. Likewise the auxiliary video burst can be approximately 60 packets in size. Utilizing these rough estimates, and ignoring the audio for simplicity, the following equation can be used to provide a starting point for tuning the queue size within the policy map:

$N * 120 \text{ packets} + M * 240 \text{ packets} = \text{Starting point for tuning the queue size for the TelePresence class.}$

Based upon this equation, a starting point for tuning the queue size within the policy-map for the TelePresence class for a circuit that supports three CTS-1000s would be:

$3 \text{ CTS-1000s} * 120 \text{ packets per CTS-1000} = 360 \text{ packets}$

As with the empirical test results for policer burst sizes shown in [Table 12-1](#), actual test results for queue-limit sizes required for support of multipoint TelePresence in a CBWFQ design have shown the equation above to be somewhat conservative. This is partly because the equation above over-simplifies the arrival and transmission of the bursts. It assumes that all packets arrive at the same time. Likewise, the equation assumes that the entire burst arrives before any of it is serialized and sent onto the circuit. In reality, as the information arrives, it is serialized. The higher the speed of the circuit, the more is serialized as it arrives, and therefore the lower the required queue depth. Finally, it must also be noted that the size of I-frames and auxiliary video bursts is variable. The values shown above represent estimates for maximum sizes seen during testing.

Despite the over-simplification of the example, the method presented can still be used as a starting point for tuning queue-limit within the policy map for supporting multipoint TelePresence. If desired the network administrator can tune the queue-limit size down to the point where TelePresence drops still do not occur.

It should be noted that increasing the queue limit size of TelePresence traffic within a CBWFQ configuration can lead to higher jitter seen by the TelePresence codecs, since a higher queue limit implies more packets can be queued, rather than discarded. At some point, jitter can lead to “late” packets which are not used by the codecs in the decoding process and may cause visible artifacts in the TelePresence video. However, decreasing the queue limit size would result in “drops” seen by the codecs, which would also result in visible artifacts in the TelePresence video. In such cases, it may be appropriate to either reconfigure TelePresence in a LLQ configuration or increase the bandwidth of the circuit such that TelePresence no longer experiences long queuing delays due to congestion when in a CBWFQ configuration.

Output Hold Queue Size Recommendations

When multipoint TelePresence is placed within a CBWFQ on a converged circuit, the output hold queue size for the WAN interface may also need to be increased above the default limit to support the bursts.

As with the queue limit size within the policy map, there is no exact science as to the size of the output hold queue to configure on the WAN interface. Tuning of the output hold queue size is dependant upon the sizes of the I-frame and auxiliary video bursts and their timing as well as the amount of congestion on the WAN circuit which causes queuing to occur. TelePresence bursts themselves may again be the cause of the temporary congestion which results in queuing.

However, in a 12-class QoS model in which TelePresence is placed in a CBWFQ, it is recommended that the output hold queue size of the WAN interface be larger than the sum of the individual queue-limit sizes of the CBWFQ classes within the policy map. This allows each of the CBWFQ classes to reach the maximum size of their queue limits and provides room for LLQ traffic without exceeding the output hold queue size of the physical WAN interface. The queue limit sizes of each of the CBWFQ classes can be determined via the show policy-map command.

TxRing Limit Recommendations

For multipoint TelePresence over a converged circuit, it is recommended to leave the TxRing configuration of the WAN interface at its default value where possible.

Multipoint TelePresence over MPLS Circuits with Ethernet Handoff

This section discusses the recommendations for support of multipoint and multiple point-to-point TelePresence over a WAN infrastructure which consists of an Ethernet handoff to a service provider MPLS network. Both dedicated and converged network configurations are discussed.

When configuring multipoint over an WAN edge design, the network administrator should follow the basic design guidance outlined in [Chapter 6, “Branch QoS Design for TelePresence.”](#) Specifically, [TelePresence Branch MPLS VPN](#) discusses specific QoS parameters as they pertain to TelePresence deployments over an MPLS VPN. The following sections simply extend those discussions to multipoint TelePresence designs.

Cisco Router and Switch Platforms Tested

The following router and switch platforms were tested for the recommendations which follow:

Cisco 7200 routers:

- Cisco 7206VXR (NPE-G2) processor (rev A) with 917504K/65536K bytes of memory
- IOS Version 12.4(4)XD7
- PA-POS-OC3MM
- PA-T3+

Cisco 3800 ISR routers:

- Cisco 3845 (revision 1.0) with 478208K bytes of memory
- IOS Version 12.4(15)T

- C3845 Mother board 1GE(TX,SFP),1GE(TX), integrated VPN and 4W Port adapter, 2 ports

Cisco 2800 ISR routers:

- Cisco 2821 (revision 53.51) with 1032192K / 16384K bytes of memory
- IOS Version 12.4(15)T
- C2821 Motherboard with 2GE and integrated VPN Port adapter, 2 ports

Multipoint TelePresence over Dedicated MPLS Circuits

When deploying multipoint TelePresence over a WAN infrastructure consisting of a dedicated MPLS circuit with an Ethernet handoff to the service provider network, egress shaping may be applied to the CE edge device. Since the physical interface rate may be a 1 Gbps Ethernet connection, but the contracted data rate to the service provider network may only be 50 Mbps, 100 Mbps, or 150 Mbps; shaping smooths out traffic to meet the contracted sub-line rate.

However, it should be noted that networking gear is not designed to buffer large amounts of traffic. Shaping should only be used to temporarily smooth bursty traffic, where the overall traffic rate is limited by mechanisms such as call admission control or manual provisioning. When shaping is active for long periods of time, the buffering capacity of routers and switches will be overrun and traffic will drop, degrading the overall TelePresence experience. In such cases, additional bandwidth should be provisioned. Note that the higher the shaped data rate, the faster the buffering capacity is exhausted on a given router or switch platform when shaping becomes active.

In a dedicated or overlay deployment, queueing is not necessarily a requirement, since all the traffic on the circuit is TelePresence related. Therefore a single-level policy map applied to the egress interface of the CE device, such as shown in [Example 12-4](#), can be implemented.

Example 12-4 Egress CE Policy Map Which Enforces Shaping Only

```
policy-map DEDICATED-EGRESS-TOP-MAP
class class-default
  shape average 50000000 1250000 6400000! All in bits/sec.
  shape max-buffers 4096
  ...
```

The network administrator needs to determine the appropriate values for the average rate (CIR), committed burst (Bc), and excess burst (Be) to best match the bursty characteristic of multipoint TelePresence. However, the network administrator is also constrained by the ingress policy map of the service provider PE device, which may police to the contracted rate. See [Example 12-5](#).

Example 12-5 Ingress PE Policy Map Which Enforces Policing Only

```
policy-map DEDICATED-TP-INGRESS
class class-default
  police cir 50000000
    conform-action set-mpls-exp-transmit 5
    exceed-action drop
  ...
```

In the policer configuration in [Example 12-5](#), the average rate (CIR) matches the contracted rate of 50 Mbps. Note that the burst parameters (Bc and Be) are not specified explicitly within the configuration and therefore assume the default value. For most Cisco IOS devices the value corresponds to a time constant (Tc) of 250 ms. For example, in the configuration above, the default committed burst (Bc) is given by the equation:

$$Bc = Tc * CIR = 250 \text{ ms} * 50 \text{ Mbps} = 12,500,000 \text{ bits} = 1,562,500 \text{ bytes}$$

Note that for policers, the burst values are typically specified in bytes not bits. The configuration also specifies that any traffic that conforms to the committed burst (Bc) is transmitted, but any traffic that exceeds it is dropped. Therefore in the example above, the excess burst (Be) is effectively not utilized.

In order to keep the PE policer from dropping ingress TelePresence traffic, and therefore stay within contract, the network design engineer needs to choose the appropriate committed rate (CIR), committed burst (Bc), and excess burst (Be) value on the shaper of the egress CE device. Two possible methods of determining these values are discussed below.

Design Option #1

With design option #1, the network administrator sets the CIR of the shaper on the CE device to match the CIR of the policer of the PE device. This is typically done when contracting a sub-line rate from a service provider.

The network administrator must ensure that the total burst (Bc + Be) sent from the CE shaper is below the committed burst (Bc) of the PE policer. However, since the time constants may be different, this may be a challenge. Configuring a long time constant on the CE shaper, perhaps 250 ms, to match the time constant of the PE policer can result in late packets seen by the remote codecs if shaping is activated at all. This nullifies the need for the shaper. In such cases it may be better not to configure a shaper at all and let the policer discard the traffic. Therefore, it may be desirable to configure a smaller time constant on the CE shaper. However a smaller time constant results in a smaller committed burst size, given by the equation $Bc = Tc / CIR$.

TelePresence bursts are relatively rare occurrences from the standpoint of the overall video which is sent at 30 frames per second. ESE testing has shown that configuring a small committed burst which yields a time constant of around 25 ms and large excess burst (Be) to handle the rest of the overall burst can work effectively in a dedicated deployment.

An example of this method using 9 CTS-1000s over a 50 Mbps MPLS circuit is discussed below. The estimated burst from 9 CTS-1000s can be estimated based on the equation:

$$9 \text{ CTS-1000s} * 128 \text{ Kbytes per CTS-1000} = 1,152,000 \text{ bytes} = 9,216,000 \text{ bits}$$

Configuring a sustained rate (CIR) of 50 Mbps with a time constant (Tc) of 25 ms yields the following committed burst (Bc) size for the CE shaper:

$$Bc = Tc * CIR = 25 \text{ ms} * 50 \text{ Mbps} = 1,250,000 \text{ bits} = 156,250 \text{ bytes}$$

Therefore, in order to hold the entire burst from 9 CTS-1000s in one time constant, the following excess burst size would need to be configured:

$$Be = 9,216,000 \text{ bits} - 1,250,000 \text{ bits} = 7,966,000 \text{ bits} = 995,750 \text{ bytes}$$

This method relies on the fact that normal TelePresence traffic is well below the recommended bandwidth allocation of 5.5 Mbps per CTS-1000 or 15 Mbps per CTS-3000 and that the bursts are typically nowhere near a full 128 Kbytes per CTS-1000 or 256 Kbytes per CTS-3000. Therefore, the service provider PE policer committed burst size (Bc) is never exceeded and TelePresence traffic is never dropped.

[Table 12-3](#) summarizes the actual test results from ESE testing of multipoint TelePresence over a dedicated 50 Mps, 100 Mbps, and 150 Mbps Ethernet handoff to an MPLS network using this methodology.

Table 12-3 ESE Tested Results for Multipoint Over Dedicated MPLS Circuits

Shaped Rate	Max CTS Table Segments	Shaper Bc (Bytes)	Shaper Be (Bytes)	Total Shaper Burst (Bytes)	Percent of Theoretical Burst Size	Policer Be (Bytes)
50 Mbps	9 CTS-1000s	156,250	800,000	956,250	83%	1,562,500
100 Mbps	3 CTS-3000s and 10 CTS-1000s	312,500	800,000	1,112,500	54%	3,125,000
150 Mbps	3 CTS-3000s and 15 CTS-1000s (See discussion)	468,750	1,400,000	1,868,750	70%	4,687,500

For all testing, the policer on the PE device was set to its default time constant of $T_c = 250$ ms and the policed average rate matched the shaped average rate of the CE device. Note that the shaper time constant was 1/10th of the policer time constant. For each test, the excess burst (Be) values were tuned down just above the point at which drops began. As can be seen, the actual burst size required ranged from 54% to 83% of the theoretical required burst sizes.

For all of the testing, the physical line rate was 1 Gbps. Because of the high physical line rate, serialization delay was low. Therefore, even though the shaped rates were 50 Mbps, 100 Mbps, and 150 Mbps, the number of CTS table segments supported was not constrained by induced jitter causing late packets on the codecs, as was the case with dedicated E3, T3, and OC-3 circuits. Instead the number of CTS table segments is bounded by the shaped bandwidth and should follow the recommendations of allocating approximately 5.5 Mbps per CTS-1000 and 15 Mbps per CTS-3000 (assuming low-speed auxiliary video input only).

Output hold queue sized for the testing were increased to their maximum values of 4,096 for all MPLS testing over a dedicate circuit with Ethernet handoff. Since TelePresence is the only traffic in a dedicated configuration, no detrimental effects are expected from this. No adjustments were made to the TxRing for the Gigabit Ethernet interface.

Design Option #2

In situations where the policer time constant (T_c) of the service provider PE has been set to prevent bursts, it may be necessary to provision a greater amount of bandwidth to support multipoint TelePresence. However, determining just how much bandwidth is required is not easy to determine. ESE has not conducted any testing with this configuration currently. The following discussion highlights some of the main issues involved.

Using the same example of 9 CTS-1000s over a 50 Mbps MPLS Ethernet handoff, if the service provider also configures a time constant (T_c) of 25 ms, then the amount of burst that can be accommodated per time constant is given by the equation:

$$B_c = T_c * CIR = 25 \text{ ms} * 50 \text{ Mbps} = 12,500,000 \text{ bits} = 156,250 \text{ bytes}$$

In this case, if the shaper were also configured with a time constant (T_c) of 25 ms, then any excess burst (Be) configured on the shaper to try to accommodate the burstiness of multipoint TelePresence would immediately be dropped by the policer. Eliminating the excess burst ($B_e = 0$) and attempting to buffer traffic in the shaper for 8 time constants (1,152,000 bytes/156,250 bytes) or approximately 200 ms would result in late packets seen by the codecs and would likely overrun the buffering capacity of the router itself.

With the ability to only handle a 156.25 Kbyte burst, the provisioning of the network is sufficient to handle a single point-to-point CTS-1000 call. It may also be able to handle two of the CTS-1000s in a multipoint call. Two CTS-1000s in a multipoint call could result in a 256 Kbyte burst when a speaker transition occurs at the same time as a PowerPoint slide transition. During the first time constant ($t = 0$ to 25 ms) the shaper would send 156.25 Kbytes and during the second time constant ($t = 25$ to 50 ms) the shaper would send the remaining 100 Kbytes. There would be additional P-slice data and audio to send as well. However, a 50 ms delay may not be sufficient to cause late packets to be seen by the codecs.

Further, by relying on the statistical assumption that speaker transitions will not likely occur at exactly the same time as PowerPoint slide transitions, the network administrator may even be able to accommodate up to four CTS-1000s in a single multipoint call. If such events do occur, some packet loss will be seen. However the network is not provisioned to handle a multipoint call with all nine CTS-1000s in a single multipoint call. Such calls will in all likelihood fail.

Since the network administrator is usually not in control of how many devices across any given circuit can be in a single multipoint call, it is recommended to design for the worse case scenario. In such cases, the only alternative may be to increase the provisioned bandwidth such that the committed burst of both the shaper and policer can accommodate the multipoint TelePresence bursts. For example, to accommodate the theoretical burst size of 1,152,000 bytes from 9 CTS-1000s in a single multipoint call, the CIR would need to be increased as follows:

$$\text{CIR} = \text{Bc} / \text{Tc} = 1,152,000 \text{ bytes} / 25 \text{ ms} = 368.64 \text{ Mbps}$$

This translates to approximately 40 Mbps per CTS-1000. At this point the design engineer may need to accept the possibility of some packet losses occurring due to a PowerPoint slide transition occurring at the same time as a speaker transition and provision the network to accommodate either a slide transition or a speaker transition occurring, but not both simultaneously. This decreases the tolerance of the network to bursts by one-half.

Going back to the example of 9 CTS-1000s, the design engineer would now design the network to support the following bursts:

$$9 \text{ CTS-1000s} * 64 \text{ Kbytes per CTS-1000} = 576,000 \text{ bytes}$$

Recalculating the amount of bandwidth required to support this configuration yields the following:

$$\text{CIR} = \text{Bc} / \text{Tc} = 576,000 \text{ bytes} / 25 \text{ ms} = 184.32 \text{ Mbps}$$

This translates to approximately 20 Mbps per CTS-1000. Finally, if this amount of bandwidth is still too high, the design engineer may redesign the network further to try and buffer the burst over two time constants, if the router has sufficient capacity to handle the burst.

Multipoint TelePresence over Converged MPLS Circuits

When deploying multipoint TelePresence over a WAN infrastructure consisting of a converged MPLS circuit with an Ethernet handoff to the service provider network, a hierarchical QoS policy may need to be applied to the CE edge device. This policy map not only shapes the overall rate of traffic to meet the contracted sub-line rate, but also polices TelePresence traffic to a particular rate within the overall contracted rate.

As previously stated in [Chapter 6, “Branch QoS Design for TelePresence,”](#) policing is not a desirable thing for TelePresence traffic. Any packets discarded by a policer degrade the TelePresence experience. However, policing may be required to limit the amount of TelePresence and/or real-time traffic sent to the service provider network when mechanisms such as CAC are not available. Also, policing may be required on certain platforms when placing traffic into an LLQ configuration, as is recommended for TelePresence.

An example of an egress hierarchical policy map applied to the interface of a CE device is shown in [Example 12-6](#).

Example 12-6 Egress CE Policy Map Which Enforces Shaping and Policing

```

policy-map EGRESS-TOP-MAP
  class class-default
    shape average 50000000 1250000 1822000
    service-policy EGRESS-CE-MODEL
policy-map EGRESS-CE-MODEL
  class VOICE
    police 1000000
    priority percent 2
  class TELEPRESENCE
    set dscp cs5
    priority percent 33          ! Places TP in LLQ
    police cir 16500000 bc 384000! Policed for 3 CTS-1000s
  class NETWORK-CONTROL
    bandwidth percent 5
  class CALL-SIGNALING
    bandwidth percent 5
  class OAM
    bandwidth percent 5
  class MULTIMEDIA-CONFERENCING
    set dscp af21
    bandwidth percent 10
  class MULTIMEDIA-STREAMING
    set dscp af21
    bandwidth percent 5
  class BROADCAST-VIDEO
    set dscp cs2
    bandwidth percent 5
  class LOW-LATENCY-DATA
    set dscp af31
    bandwidth percent 5
  class HIGH-THROUGHPUT-DATA
    bandwidth percent 4
  class SCAVENGER
    bandwidth percent 1
  class class-default
    bandwidth percent 20
    random-detect
  ...

```

The policy map is similar to that which would be applied to support a single TelePresence device in a point-to-point configuration, however it has been scaled to support multipoint TelePresence.

Policer Details

The average rate (CIR) and committed burst (Bc) of the policer within the TelePresence traffic class have been scaled for support of multipoint TelePresence. Scaling of the average rate can be accomplished by utilizing 5.5 Mbps per CTS-1000 and 15 Mbps per CTS-3000. Keeping in mind that it is recommended to limit real-time traffic to approximately 1/3 of the bandwidth of any link, the configuration example above is designed to support a maximum of 3 CTS-1000s * 5.5 Mbps per CTS-1000 = 16.5 Mbps in the TelePresence traffic class. This is approximately 33% of the overall 50 Mbps shaped rate. In this example, if more than 3 CTS-1000s need to be supported in an LLQ configuration, a higher contracted rate to the service provider would be recommended.

Scaling the committed burst (Bc) of the policer within the TelePresence traffic class can also be accomplished by utilizing the burst estimates of 128 Kbytes per CTS-1000 and 256 Kbytes per CTS-3000. In the example above, which supports 3 CTS-1000s, this translates to 3 CTS-1000s * 128 Kbytes per CTS-1000 = 384 Kbytes.

Shaper Details

The configuration of the shaper follows the guidance presented in [Design Option #1](#) for support of multipoint TelePresence over a dedicated MPLS link. The average rate (CIR) is set to the overall contracted rate from the service provider. The time constant (Tc) is selected to be relatively low, 25 ms in the example above. This results in a relatively small committed burst (Bc) of 1,250,000 bits or 156,250 bytes. In order to accommodate bursts from the policer which has been configured for 384,000 bytes, an excess burst of $384,000 - 156,250 = 227,750$ bytes or 1,822,000 bits is configured.

Tested Results

In reality, both the shaper excess burst (Be) and the policer committed burst (Bc) may be tuned down somewhat. [Table 12-4](#) summarizes the actual test results from ESE testing of multipoint TelePresence over a shared 50 Mps, 100 Mbps, and 150 Mbps Ethernet handoff to an MPLS network using the methodology discussed above.

Table 12-4 ESE Tested Results for Multipoint over Shared MPLS Circuits

Shaper CIR	Shaper Bc	Shaper Be	Total Shaper Burst	Policer CIR	Policer Bc	Maximum CTS Table Segments	Output Hold Queue
50 Mbps	156.25 Kbytes	156.25 Kbytes	312.5 Kbytes	16.5 Mbps	192 Kbytes	3 CTS-1000s	40
100 Mbps	312.5 Kbytes	365.5 Kbytes	678 Kbytes	33 Mbps	384 Kbytes	6 CTS-1000s	80
150 Mbps	467.75 Kbytes	800 Kbytes (See discussion)	1,268.75 Kbytes	49.5 Mbps	805 Kbytes	9 CTS-1000s	4096 (See discussion)

As can be seen from [Table 12-4](#), testing has indicated that the policer committed burst (Bc) can normally be tuned below the 128 Kbyte per CTS-1000 recommendation. However, the shaper excess burst has to also account for the other CBWFQ classes of traffic, Voice LLQ class traffic, as well as the TelePresence LLQ class traffic. Therefore the overall burst capability of the shaper is above the policer burst capability in a converged circuit configuration.



Note

The Be value of 800 Kbytes for the 150 Mbps test was not tuned down to a minimum. It is believed that a value of 436 Kbytes would have worked as well.

Output Hold Queue and TxRing Tuning

In a converged configuration, the output hold queue will likely need to be tuned above the default of 40 packets for a Gigabit Ethernet interface. [Table 12-4](#) presents the output hold queue sizes required for the 50 and 100 Mbps tests. However, the results are highly dependant on the overall traffic on the interface. The output hold queue for the 150 Mbps test was simply increased to its maximum value of 4,096 to see if there were any detrimental effects on TelePresence quality by simply increasing the output hold queue to its maximum value. No adverse effects were observed. However, since TelePresence was placed in the LLQ, no adverse effects from queueing delays were expected either.

Platform and Linecard Test Results and Recommendations

This section summarizes the recommendations for support of multipoint TelePresence over various switch linecards and platforms within the LAN; as well as various switch and router platforms within the WAN. All recommendations apply to converged network designs.



Note

The platforms and linecards listed in the tables below have been tested by ESE. Just because a platform, linecard, or particular WAN media has not been tested, does not necessarily mean that it will not work in a multipoint TelePresence design.

Multipoint TelePresence over Switch Linecards and Platforms within the LAN

The reader is encouraged to review [Chapter 5, “Campus QoS Design for TelePresence”](#) for a detailed discussion of the queuing structure and per-port buffering of each switch platform and/or linecard, as well as detailed recommendations for configuring each switch platform and/or linecard to support TelePresence. The configuration recommendations presented in [Chapter 5, “Campus QoS Design for TelePresence”](#) apply to multipoint TelePresence as well.

[Table 12-5](#) and [Table 12-6](#) summarize the results from the platform and linecard testing for support of multipoint TelePresence, both for use at the LAN head-end (connected to the CTMS) and for the LAN remote-site (connected to CTS endpoints).



Note

ESE testing of all linecards and/or platforms at the head-end was performed up to 24 table segments or one-half the current capacity of a fully loaded CTMS.

Table 12-5 *Modular Switch Platforms*

Platform	Supervisor	Linecard	Connected to CTMS	Connected to CTS Endpoints
Catalyst 6500	Sup-32	WS-X6148A-GE-TX	Tested and Caveats (See discussion)	Tested and Recommended
	Sup-32	WS-X6548-GE-TX	Tested and Caveats (See discussion)	Tested and Recommended
	Sup-720	WS-X6748-GE-TX	Tested and Recommended	Tested and Recommended
Catalyst 4500	SupII+	WS-X4548-GB-RJ45V	Tested and Caveats (See discussion)	Tested and Recommended
	SupII+	WS-X4448-GB-RJ45	Tested and Caveats (See discussion)	Recommended

Table 12-6 Fixed Configuration Switch Platforms

Platform	Connected to CTMS	Connected to CTS Endpoints
Catalyst 4948-10GE	Tested and Recommended	Tested and Recommended
Catalyst 3750G-48PS-E	Tested and Caveats (See discussion)	Tested and Caveats (See discussion)

Both the WS-X6748-GE-TX linecard of the Catalyst 6500 platform and the Catalyst 4948-10GE platform are recommended to connect either a CTMS at the head-end of a multipoint TelePresence deployment or multiple CTS endpoints at the remote-end of a multipoint deployment, in a converged network design.

The WS-X6548-GE-TX and WS-X6148A-GE-TX linecards of the Catalyst 6500 platform are recommended to connect multiple CTS endpoints at the remote-end of a multipoint deployment in a converged network design. However, since these linecards are 8:1 oversubscribed, they are not recommended to connect a CTMS at the head-end of a multipoint deployment unless the network administrator dedicates the entire range of eight ports corresponding to the 1 Gbps uplink to the switch fabric shared by these eight ports. Port ranges are 1-8, 9-16, 17-24, 18-32, 33-40, or 41-48.

The WS-X4548-GB-RJ45V and WS-X4448-GB-RJ48 linecards of the Catalyst 4500 platform are recommended to connect multiple CTS endpoints at the remote-end of a multipoint deployment in a converged network design. However, since these linecards are again 8:1 oversubscribed, they are not recommended to connect a CTMS at the head-end of a multipoint deployment unless the network administrator dedicates the entire range of eight ports corresponding to the 1 Gbps uplink to the switch fabric shared by these eight ports. Port ranges are 1-8, 9-16, 17-24, 18-32, 33-40, or 41-48.

The Catalyst 3750G-48PS-E is recommended to connect either a CTMS at the head-end of a multipoint TelePresence deployment or multiple CTS endpoints at the remote-end of a multipoint deployment, in a converged network design. However, testing was only performed in a non-stacked configuration. Due to potential issues regarding oversubscription of the dual 16 Gbps ring backplane, it is not recommended to use the Catalyst 3750G-48PS-E in a stacked configuration either for connectivity to the CTMS or CTS endpoints.

Table 12-7 summarize the results from supervisor and linecard ports used as uplinks between switches when supporting multipoint TelePresence.

Table 12-7 Uplink Ports between Switches

Supervisor or Linecard Port	Results
Catalyst 6500 Sup-32 10GE Port	Tested and Recommended
Catalyst 6500 Sup-32 1 GE Port	Tested and Caveats (See discussion)
WS-X6704-10GE	Tested and Recommended
WS-X6708-10GE	Tested and Recommended
Catalyst 6500 Sup-720 GE Port	Tested and Caveats (See discussion)
WS-X6748-GE-TX	Tested and Caveats (See discussion)
Catalyst 4500 SupII+ 10 GE Port	Tested and Recommended
Catalyst 4948-10GE Platform 10 GE Port	Tested and Recommended
Catalyst 3750G-48PS-E Platform GE Port	Tested and Caveats (See discussion)

The 10 GE uplink ports on the Catalyst 6500 Sup-32 and Catalyst 4500 SupII+ are both recommended for uplinks between switches. Likewise ports on the Catalyst 6500 WS-X6704-10GE and WS-X6708-10GE linecards are recommended for use as uplink ports between switches. Finally the 10 GE uplink ports on the Catalyst 4948-10GE platform are also recommended as uplinks between switches.

The 1 GE uplink ports on the Catalyst 6500 Sup-32 and Sup-720, 1 GE ports on the WS-X6748-GE-TX, and 1 GE ports on the Catalyst 3750G-48PS-E platform can also be used as uplink ports between switches. However, caution should be applied when utilizing 1 gigabit Ethernet ports for connectivity to endpoint devices (including TelePresence CTS units) and only providing a 1 gigabit Ethernet uplink port between switches. The amount of oversubscription of the uplink port can be high, resulting in congestion. Multiple gigabit Ethernet ports can be implemented in a Gigabit EtherChannel configuration with some caveats.

[Example 12-7](#) shows part of the configuration of a Catalyst 3750G switch for a Gigabit EtherChannel configuration consisting of three ports.

Example 12-7 Catalyst 3750G Gigabit EtherChannel Configuration

```

!
port-channel load-balance src-dst-ip
!
!
interface Port-channel1
 no switchport
 ip address 10.16.3.1 255.255.255.0
 load-interval 30
!
!
interface GigabitEthernet1/0/33
 no switchport
 no ip address
 load-interval 30
 shutdown
 srr-queue bandwidth share 1 30 35 5
 priority-queue out
 mls qos trust dscp
 channel-group 1 mode auto
!
interface GigabitEthernet1/0/34
 no switchport
 no ip address
 load-interval 30
 shutdown
 srr-queue bandwidth share 1 30 35 5
 priority-queue out
 mls qos trust dscp
 channel-group 1 mode auto
!
interface GigabitEthernet1/0/35
 no switchport
 no ip address
 load-interval 30
 srr-queue bandwidth share 1 30 35 5
 priority-queue out
 mls qos trust dscp
 channel-group 1 mode auto
!

```

Depending upon the platform, load balancing on the port-channel group can be done in various ways—by source IP address, by destination IP address, by source MAC address, by destination MAC address, by source and destination IP address, or by source and destination MAC address.

There are two main issues with regard to the use of EtherChannel technology in this configuration. First, EtherChannel technology does not take into account the bandwidth of each flow. Instead, it relies on the statistical probability that, given a large number of flows of relatively equal bandwidths, the load is equally distributed across the links of the port-channel group. Since the CTMS represents a single IP address and a single MAC address, load balancing traffic based on source IP address or source MAC address on the port channel is not considered to be effective. All outbound traffic from the CTMS would take a single path. This could result in a lopsided distribution of traffic, in which one link could have more than 1 Gbps of traffic, since video flows are variable in nature, while other links have relatively little traffic. Depending upon the configuration of the network past the EtherChannel group, the destination MAC address in all likelihood corresponds to a router port. Therefore, load-balancing based on the destination MAC address or source and destination MAC address is also not considered to be effective. Load balancing based on the destination IP address, or source and destination IP address, are considered to be effective for this type of design, since the CTMS sends traffic to multiple CTS endpoint IP addresses. Finally, it should be noted that each source and destination IP address represents a video connection between the CTMS and one CTS endpoint. Therefore there is no chance of out-of-sequence packets being seen by the CTMS or CTS units in this configuration during normal operation.

The second issue with EtherChannel technology is that it does not take into account any QoS configuration on the individual Gigabit Ethernet links. Again, it relies on the statistical probability that, given a large number of flows of with different QoS markings, the load of those individual flows is equally distributed across the links of the port-channel group. Given a failover situation in which one of the links of an EtherChannel group fails, the sessions crossing that link would be re-allocated across the remaining links. Since EtherChannel technology has no awareness of QoS markings, it could easily re-allocate more real-time flows across any one of the links than the link is configured to accommodate. This would result in degraded real-time services such as TelePresence and voice, although there is sufficient real-time bandwidth within the EtherChannel group to easily accommodate the real-time traffic. This can happen in a non-failover situation as well. Therefore, caution should be used when deciding to utilize EtherChannel technology versus a single higher-speed uplink port.

Multipoint TelePresence over WAN Switch and Router Platforms

Table 12-8 summarizes the recommendations for WAN router and switch platforms along with associated media tested for support of multipoint TelePresence in a converged network design.

Table 12-8 Platforms and WAN Media

Platform	WAN Hardware	WAN Media	Results
Cisco 2821 ISR	Integrated Gigabit Ethernet Interface	MPLS Handoff with 50 Mbps Shaped Fast Ethernet	Tested and Caveats (See discussion)
Cisco 3845 ISR	Integrated Gigabit Ethernet Interface	MPLS Handoff with 50 Mbps Shaped Gigabit Ethernet	Tested and Recommended
Cisco 3845 ISR	Integrated Gigabit Ethernet Interface	MPLS Handoff with 100 Mbps Shaped Gigabit Ethernet	Tested and Recommended
Cisco 3845 ISR	Integrated Gigabit Ethernet Interface	MPLS Handoff with 150 Mbps Shaped Gigabit Ethernet	Tested and Recommended (See Discussion)

Table 12-8 Platforms and WAN Media

Platform	WAN Hardware	WAN Media	Results
Cisco 7206VXR with NPE-G2	Integrated Gigabit Ethernet Interface	MPLS Handoff with 50 Mbps Shaped Gigabit Ethernet	Tested and Recommended
Cisco 7206VXR with NPE-G2	Integrated Gigabit Ethernet Interface	MPLS Handoff with 100 Mbps Shaped Gigabit Ethernet	Tested and Recommended
Cisco 7206VXR with NPE-G2	Integrated Gigabit Ethernet Interface	MPLS Handoff with 150 Mbps Shaped Gigabit Ethernet	Tested and Recommended
Cisco 7206VXR with NPE-G2	PA-T3+	Clear Channel T3	Tested and Recommended
Cisco 7206VXR with NPE-G3	PA-POS-OC3MM	OC-3 POS	Tested and Recommended
Catalyst 6506-E with Sup-32	7600-SIP-200 with SPA-4XT3/E3	Clear Channel E3	Tested and Recommended
Catalyst 6506-E with Sup-32	7600-SIP-400 with SPA-2XOC3-POS	OC-3 POS	Tested and Recommended
Catalyst 6506-E with Sup-32	7600-SIP-400 with SPA-1XOC12-POS	OC-12 POS	Tested and Recommended
Catalyst 6506-E with Sup-32	7600-SIP-400 with SPA-1XOC48-POS	OC-48 POS	Tested and Recommended

For all tests involving handoff to an MPLS network via either a 50 Mbps, 100 Mbps, or 150 Mbps Ethernet link, an integrated Gigabit Ethernet interface on the Cisco 2821 ISR, Cisco 3845 ISR, or Cisco 7206VXR was utilized.

Significant input errors (overruns) were observed on the interface of the Cisco 2821 ISR at 1 Gbps physical line speeds during the testing when both background traffic and TelePresence traffic were applied. Based on these results, it is not recommended to utilize the Cisco 2821 ISR for a 50 Mbps shaped Ethernet handoff to a service provider MPLS network when the service provider physical line rate is 1 Gbps in a converged network design. Instead it is recommended to use the Cisco 2821 ISR for a 50 Mbps shaped Ethernet handoff to the service provider MPLS network only when the service provider physical line rate is 100 Mbps. Due to these overrun errors, the Cisco 2821 is also not recommended for Ethernet handoff to an MPLS network above a 50 Mbps shaped rate in a converged network design.

The Cisco 3845 ISR experienced a small number of input errors (overruns) during testing at 150 Mbps shaped Ethernet handoff to the MPLS network. However, the input errors were well below the 0.1% threshold for the drop rate of TelePresence. Therefore the Cisco 3845 ISR is recommended for shaped 50 Mbps, 100 Mbps, and 150 Mbps Ethernet MPLS handoff to a service provider network.

The Cisco 7206VXR experienced no input errors during any of the tests with a shaped Ethernet handoff to the MPLS network. Therefore the Cisco 7206VXR is recommended for shaped 50 Mbps, 100 Mbps, and 150 Mbps Ethernet MPLS handoff to a service provider network as well.



CHAPTER 13

Internal Firewall Deployments with Cisco TelePresence

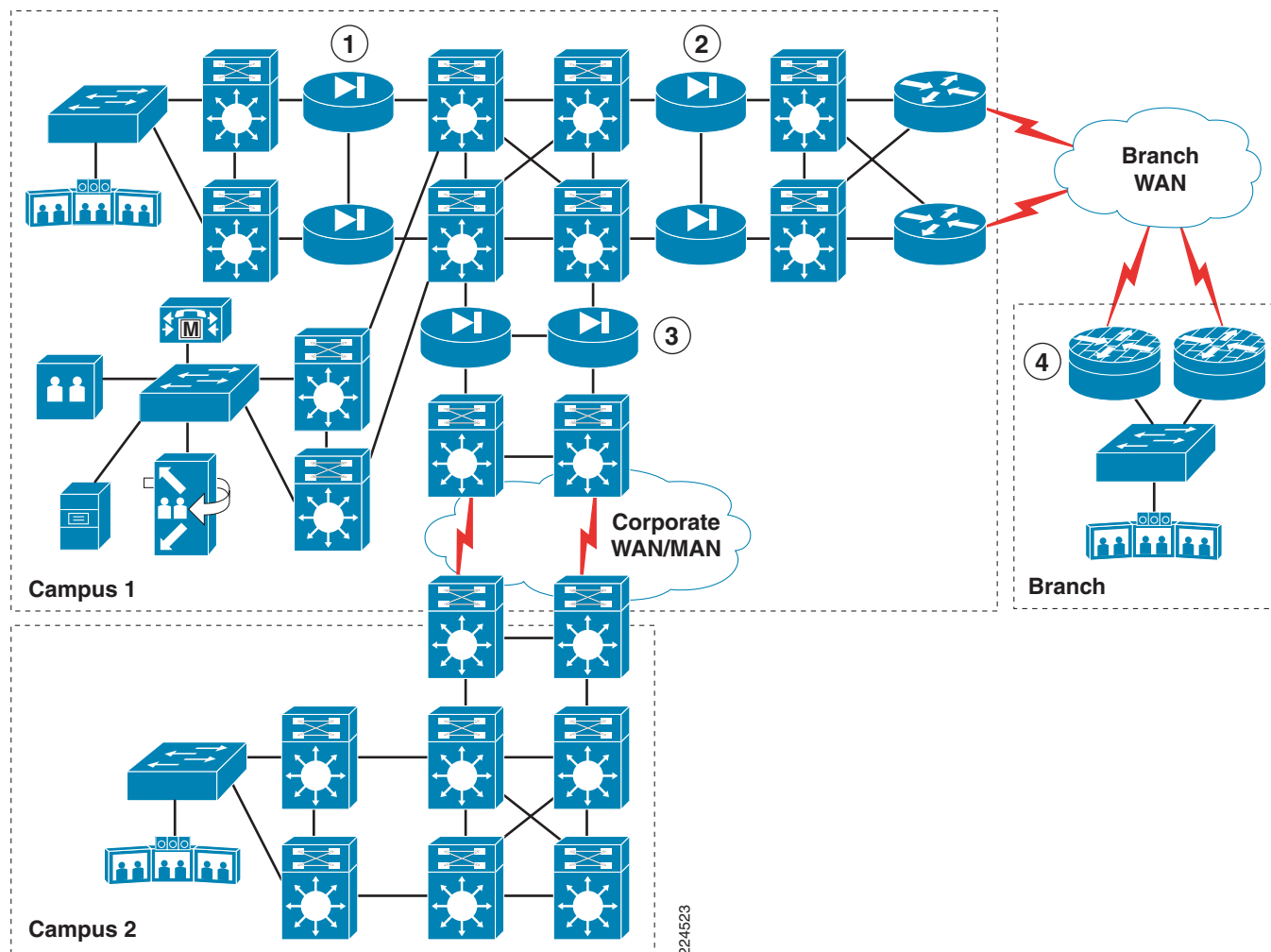
Overview

This chapter t discusses the deployment of Cisco TelePresence within an enterprise organization which may implement internal firewalling. Internal firewalling may be deployed within an organization for a number of reasons, including:

- Access control within an enterprise campus
Firewalling may be implemented within an enterprise campus in order to provide access control for a corporate department, division, or service module within the campus network. This may be done for regulatory or internal security reasons.
- Access control between enterprise campus locations
Firewalls may be implemented to provide NAT services between two campus locations. For example, when two companies merge, there may be a period of time where NAT is done due to an overlapping IP address space between the two sides of the company. Access control between the two campus locations may also be implemented if necessary.
- Access control from branch locations to corporate campus sites
Firewalling may be implemented at WAN aggregation points within an enterprise organization in order to restrict inbound access from branch locations to certain protocols and/or resources within the corporate campus locations. This may be done to enhance the trust boundary between the branch locations and the corporate campus.
- Access control within enterprise branch locations
Firewalling may be implemented within branch locations in order to restrict access to certain devices within a branch. An example would be the isolation of IP-based point-of-sale terminals within a store location.

[Figure 13-1](#) shows the deployment of firewalling in these four areas in relation to a TelePresence deployment.

Figure 13-1 Internal Firewalling within an Enterprise Organization



The following flows may need to be enabled through the firewall in order to support the TelePresence deployment:

- Device provisioning flows between the Cisco Unified Communications Manager (CUCM) cluster and CTS endpoints in order to successfully bring the CTS endpoints onto the network and register with CUCM.
- Call scheduling and service flows between the CTS Manager and the CTS endpoints.
- Call signaling flows between CTS endpoints and the CUCM cluster in order to initiate and terminate TelePresence meetings.
- The actual audio and video media flows between CTS endpoints in a point-to-point call.
- The actual audio and video media flows between CTS endpoints and the Cisco TelePresence Multipoint Switch (CTMS) in a multipoint call.
- Flows between network management stations and the CTS endpoints to successfully manage the TelePresence deployment.

Cisco Firewall Platforms

Cisco currently provides three firewall product lines:

- ASA 5500 Series of firewall appliances
- IOS Firewall running on Cisco IOS router platforms
- Firewall Service Module (FWSM) for the Catalyst 6500 switches and Cisco 7600 Series routers

This chapter includes test results from the ASA 5500 Series of firewall appliances only. The recommendations do not apply to the FWSM or IOS Firewall, which have not been tested within the ESE lab.

The following firewall platform and software version was tested for the recommendations within this document:

- ASA5550, 4096 MB RAM, CPU Pentium 4 3000 MHz
- Software Version ASA722-K8

**Note**

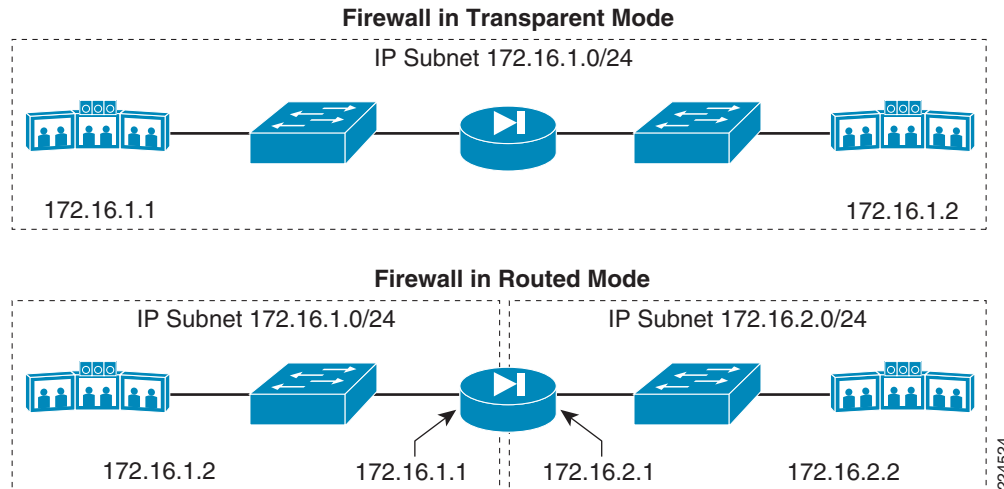
The Cisco PIX firewall series of appliances have not been tested since they have been replaced with the ASA 5500 Series of firewall appliances.

Firewall Deployment Options

The following sections discuss some of the options available when deploying a firewall within an Enterprise. Where appropriate, information regarding the affect of the choice of firewall deployment on TelePresence deployments is discussed.

Transparent versus Routed Mode

Firewalls such as the Cisco ASA 5500 Series can operate either as a Layer 2 (transparent mode) or a Layer 3 (routed mode) device. A firewall operating in transparent mode does not mean that access control decisions are made based upon Layer 2 MAC address information. Access control decisions are still made based upon Layer 3 and higher (in the case of Application Layer Protocol Inspection) information. A firewall operating in transparent mode has the same IP subnet on both sides of the firewall as shown in [Figure 13-2](#). This is often beneficial for deploying firewalling in existing networks, since IP addressing does not have to be changed to insert the firewall.

Figure 13-2 Transparent Mode versus Routed Mode Firewall

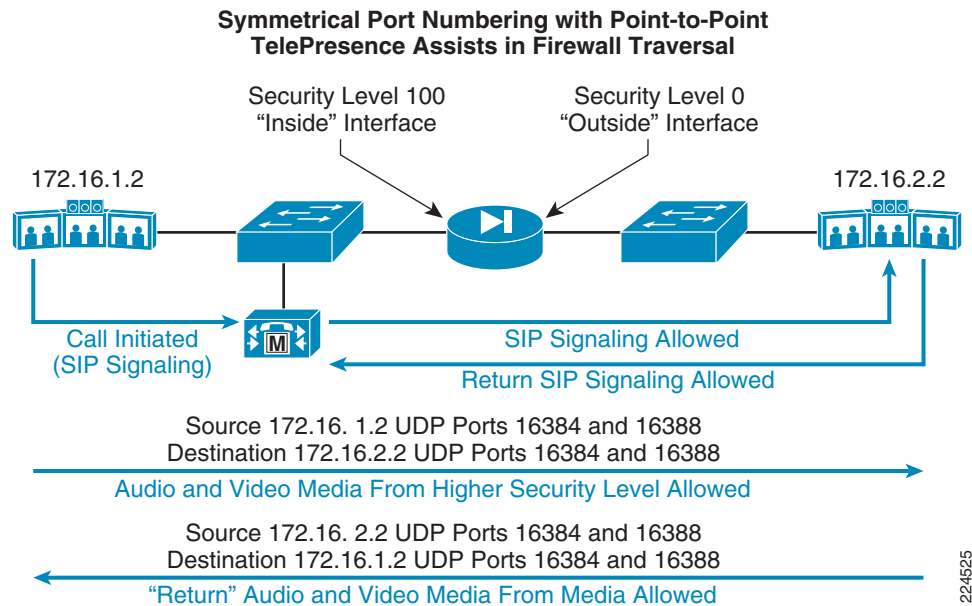
A firewall operating in routed mode has different IP subnets on both sides of the firewall, as shown in [Figure 13-2](#). This is the more traditional mode of operation of a firewall and was the only configuration tested and discussed within this document.

Equal versus Unequal Interface Security Levels

Firewalls such as the Cisco ASA 5500 Series operate based upon the concept of levels of trust within the network. This is reflected through security levels configured on firewall interfaces. Security levels range from 100, which is the most secure interface, to 0, which is the least secure interface. These are often referred to as the “inside” and “outside” interfaces on a firewall with only two interfaces.

By default, traffic initiated from a device on an interface with a higher security level is allowed to pass to a device on an interface with a lower security level. Return traffic corresponding to that session is dynamically allowed from the lower interface security level to the interface with the higher security level. Note that the term “session” can also apply to return UDP-based traffic, although there is no actual session as with TCP-based traffic. The use of symmetric port numbering in point-to-point TelePresence calls can actually provide a benefit when traversing firewalls due to this behavior. This is shown in [Figure 13-3](#).

Figure 13-3 Symmetrical TelePresence Ports and Firewalls



In [Figure 13-3](#), since the UDP ports are symmetric, the video generated by the CTS endpoint on the interface with the lower security level appears as if it is the return traffic of the CTS endpoint on the higher security level and is therefore allowed back through the firewall. However, it should be noted that this symmetric use of ports does not necessarily hold for multipoint TelePresence calls.

By default, traffic initiated from a device on an interface with a lower security level is not allowed to pass to a device on an interface with a higher security level. This behavior can be modified with an ingress access-control list (ACL) on the lower security interface level. For example, in [Figure 13-3](#), an ingress ACL on the interface with the lower security level is needed to allow both the primary codec and the associated IP phone of the CTS endpoint to register with the CUCM cluster via the SIP protocol. As well as an inbound ACL, static translations may need to be configured within the firewall to allow the devices on the "inside" network to be visible to the "outside" network.

Depending upon the firewall deployment, an ingress ACL applied to the interface with the higher security level may also be desired to limit traffic going from higher level security interfaces to interfaces with lower security levels.

Cisco ASA 5500 Series firewalls can also be allowed to operate with interfaces having equal security levels. By default traffic is not allowed to pass between interfaces having the same security level unless the same-security-traffic permit inter-interface global command is also configured within the firewall. Again, ingress ACLs applied on each interface and static translations can be used to specifically allow access between certain devices and protocols connected to interfaces with equal security levels.

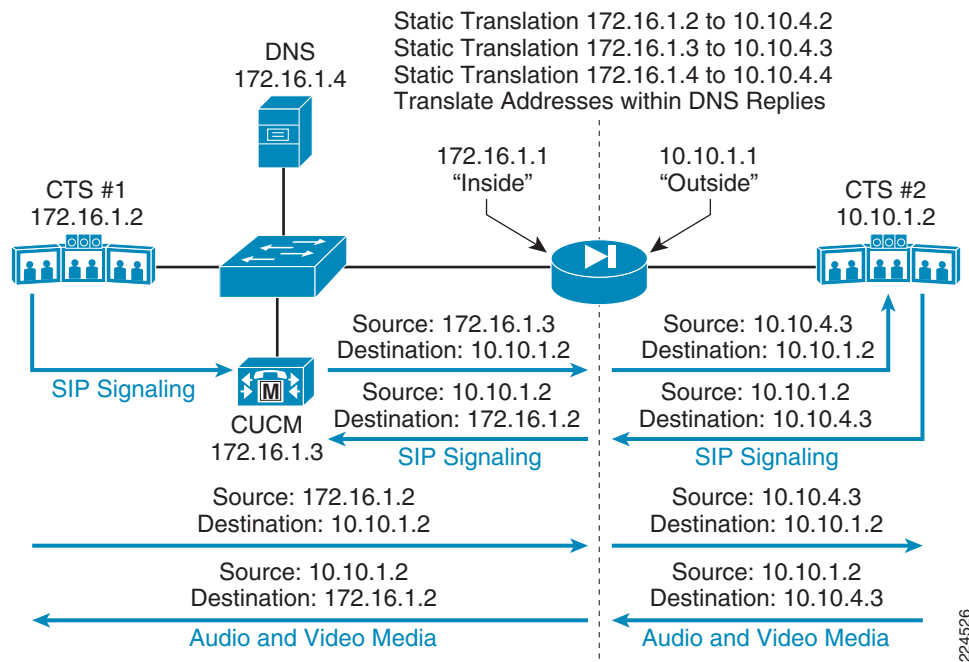
TelePresence has been tested in point-to-point configurations using both unequal and unequal cost interface configurations on the ASA 5000 firewall. Specific protocols which may be required within ACLs are discussed in [TelePresence Protocol Requirements](#).

Network Address Translation (NAT)

NAT is normally used to hide the internal IP addressing of an enterprise organization when accessing resources on the Internet. It is not utilized often in internal firewall deployments, unless overlapping IP address ranges exist within the enterprise. These could temporarily result from acquisitions and mergers between companies.

NAT utilizes the concept of dynamic address pools which are used to translate “local” addresses to “global” addresses, as well as individual static translations between devices on the “local” side and the “global” side. [Figure 13-4](#) provides an example of one-sided NAT within a TelePresence deployment.

Figure 13-4 TelePresence Deployment with NAT



Note that in [Figure 13-4](#), the access-control lists allowing appropriate inbound and outbound addresses and protocols are not shown in order to simplify the figure.

When deploying TelePresence within a one-sided NAT environment, static translations are needed between the “inside” or “local” IP addresses of the CTS endpoint, associated IP 7970 Phone, the CUCM, and possibly the DNS server. These translate to the “global” IP addresses which are visible to the network on the “outside” interface. In this type of deployment, the “inside” network is aware of the entire IP address range of the “outside” network and proper IP routing must be configured for reachability. The “outside” network, however, is not aware of the “inside” IP addressing. All “inside” addresses are translated to “global” IP addresses by the firewall. Therefore, IP routing on the outside must only be able to reach the IP subnet of the “global” address pool. Note that many-to-one NAT, also referred to as Port Address Translation (PAT), does not work in this environment due to the requirements for static translations.

The ASA 5500 firewall also has an option for translating IP addresses within DNS replies. This allows the DNS server to be deployed on the “inside” of the NAT firewall and still correctly hand out IP addresses to devices “outside” the firewall. Modifications to the configuration of the CTS endpoint are also required for this configuration to operate correctly. The DNS server entry in CTS#2 needs to be configured to point to the “global” address of the DNS server (10.10.4.4 in the example above). The entry for the CUCM cluster configured in CTS#2 can utilize the hostname of the CUCM server. The DNS translation function of the firewall will ensure that the “local” CUCM address is translated to the “global” address before it is handed to CTS#2. Otherwise, the configuration of CTS#2 would have to be modified to include the “global” address of the CUCM server as well. An alternative to translation of DNS replies within the firewall is a split DNS implementation. In this configuration a DNS server is deployed on both sides of the firewall with the appropriate records for the particular addressing used on that side of the firewall. However, this method doubles the amount of administrative work in maintaining DNS across the enterprise.

Point-to-Point TelePresence has been tested in one-sided NAT configurations as well as non-NAT configurations.

Application Layer Protocol Inspection

Application layer protocol inspection is required for services that embed IP addressing information in user data sections of a packet or that open secondary channels on dynamically assigned ports. These protocols require the firewall to perform a deep packet inspection in order to extract such information. The Cisco TelePresence solution embeds IP addressing information within the SIP signaling between CTS endpoints and the CUCM cluster. SIP signaling is also used to dynamically specify the RTP audio and video media ports, which range from UDP port 16384 to 32677. The firewall therefore dynamically opens and closes the required audio and video media ports based upon inspection of the SIP signaling between the CTS endpoints and the CUCM cluster.

Without SIP inspection, the network administrator may have to statically open these UDP port ranges for the IP addresses corresponding to each CTS endpoint which needs to pass through the firewall. This presents a larger security vulnerability from a firewall perspective. Therefore, when possible, the use of application layer protocol inspection of SIP traffic for TelePresence deployments is recommended. Point-to-point TelePresence has been tested in configurations with and without application layer protocol inspection of SIP traffic.

Note also that DNS and TFTP application layer protocol inspection are enabled by default with the ASA 5500 Series of firewall appliances. It is recommended to leave these enabled.

TelePresence Protocol Requirements

The following sections discuss some of the protocols which may need to be enabled through a firewall, depending upon the configuration of the network and CTS endpoints. The discussion is geared toward identification of TelePresence protocol requirements based upon the following functions:

- Provisioning the devices onto the network
- Registering the devices with the CUCM
- Call scheduling and services flows
- Call signaling flows needed to initiate and terminate a TelePresence meeting
- Support of the actual media flows across the network
- Management of the devices on the network

Device Provisioning Flows

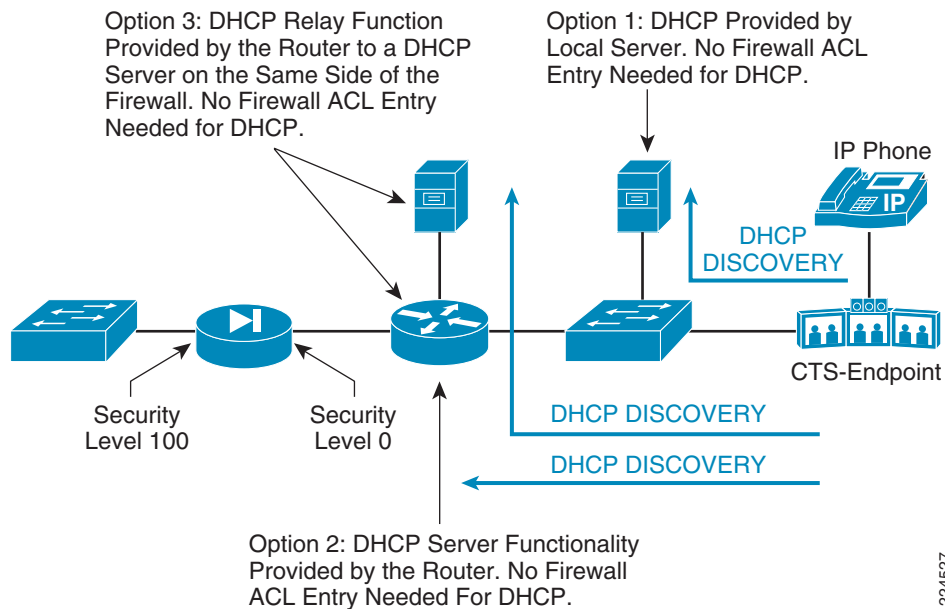
Device provisioning flows are the signaling and data flows needed by CTS endpoints to boot up and register with the CUCM cluster.

Dynamic Host Configuration Protocol (DHCP)

DHCP is used for dynamic IP address assignment. Client-sent DHCP packets have UDP source port 68 and destination port 67. Server-sent DHCP packets have UDP source port 67 and destination port 68. If both the primary codec of the CTS-endpoint as well as its associated IP phone use static IP addressing, then DHCP is not required. If either the primary codec or its associate IP phone use dynamic IP

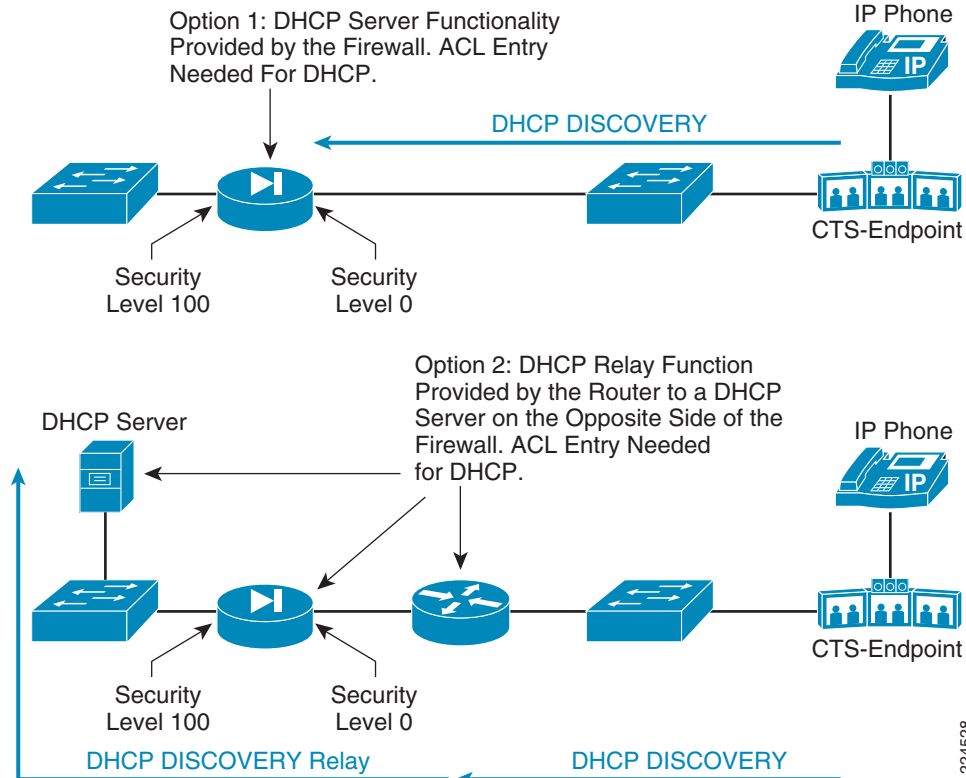
addressing, then DHCP may be required to pass through the firewall, depending upon the network configuration. [Figure 13-5](#) shows three DHCP deployment options where no firewall ACL entries are needed to support TelePresence.

Figure 13-5 DHCP Deployments Not Requiring Firewall ACL Entries for TelePresence



The initial DHCP DISCOVERY packet is sent by the CTS endpoint to the IP broadcast address (255.255.255.255). Therefore either the DHCP server functionality must be local to the IP subnet or a router can be configured to provide DHCP relay functionality to a DHCP server on another IP subnet. As long as the DHCP server functionality is on the same side of the firewall, no ACL entries are needed for DHCP flows from TelePresence devices.

[Figure 13-6](#) shows two DHCP deployment options where firewall ACL entries are needed to support TelePresence.

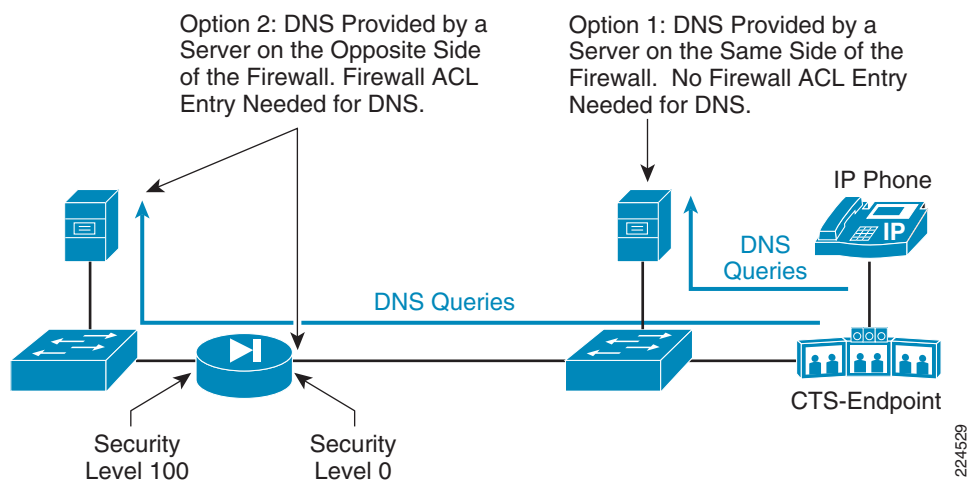
Figure 13-6 DHCP Deployments Requiring Firewall ACL Entries for TelePresence

Option 1 of [Figure 13-6](#) shows that even if the firewall provides the DHCP server functionality, an ACL entry may be needed to allow the inbound DHCP DISCOVERY packet as well as further DHCP packet exchanges. If a router provides DHCP relay functionality to a DHCP server on the opposite side of a firewall, as shown in Option 2, an ACL entry may be needed to allow the relayed DHCP exchanged to occur as well.

Domain Name System (DNS)

DNS utilizes UDP port 53 for hostname to IP address resolution. When DHCP returns an IP address to a CTS endpoint, it can also provide further information, such as the CUCM cluster to which the CTS-endpoint must register and its configuration download server (typically the CUCM cluster again). This information can be provided in the form of IP addresses or hostnames. Note that DHCP typically also provides the IP address of the DNS server and the domain name of the CTS endpoint. Alternatively, the network administrator can statically configure either the IP address or hostnames of the CUCM cluster and configuration download server within the CTS endpoints. Regardless of which method used, if hostnames are provided to the CTS endpoint, DNS is required in order for the CTS endpoint to resolve the hostname to an IP address.

DHCP may be required to pass through the firewall, depending upon the network configuration. DNS queries are initiated by the CTS endpoints. Therefore, if the DNS server is located on the same side of the firewall as the CTS endpoint, no ACL entry may be needed for DNS. However, if the DNS server is located on the opposite side of the firewall from the CTS endpoint, an ACL entry allowing DNS queries from the CTS endpoint may be needed. These two deployment options are shown in [Figure 13-7](#).

Figure 13-7 DNS support for TelePresence in a Firewall Deployment**Note**

As mentioned previously, application layer protocol inspection of DNS packets is on by default within the ASA 5500 Series of firewall appliances.

Configuration Download Protocols

Configuration download protocols include those needed by the CTS codecs as well as their attached IP phones to upgrade system load images and download device configuration files. The CUCM cluster often provides the download server functionality and is therefore the only example discussed within this chapter. Further, this section of the discussion assumes that the CUCM cluster is on the opposite side of the firewall from the CTS endpoint.

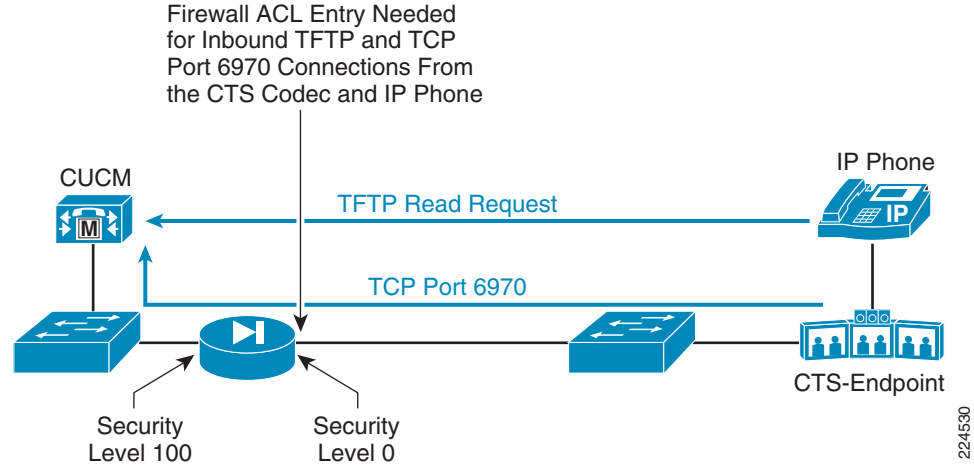
Trivial File Transfer Protocol (TFTP)

The IP 7970 Phone which serves as the user interface for TelePresence calls requires the TFTP protocol for downloading system images during upgrades as well as for downloading configuration files. TFTP servers listen on UDP port 69, but then dynamically assign a different port number for actual file transfers. Therefore, it is recommended to leave application layer protocol inspection enabled for TFTP within the firewall, which is the default for ASA 5500 Series firewall appliances. However, an ACL entry allowing the initial TFTP read request from the CTS endpoint to the CUCM cluster may be needed.

TCP Port 6970

Unlike IP phones, TelePresence codecs do not use TFTP for download of system images and configuration files. TelePresence codecs utilize HTTP over TCP port 6970 to download system images and configuration files. An ACL entry allowing the session to be initiated by the CTS endpoint on the opposite side of the firewall from the CUCM cluster may be required.

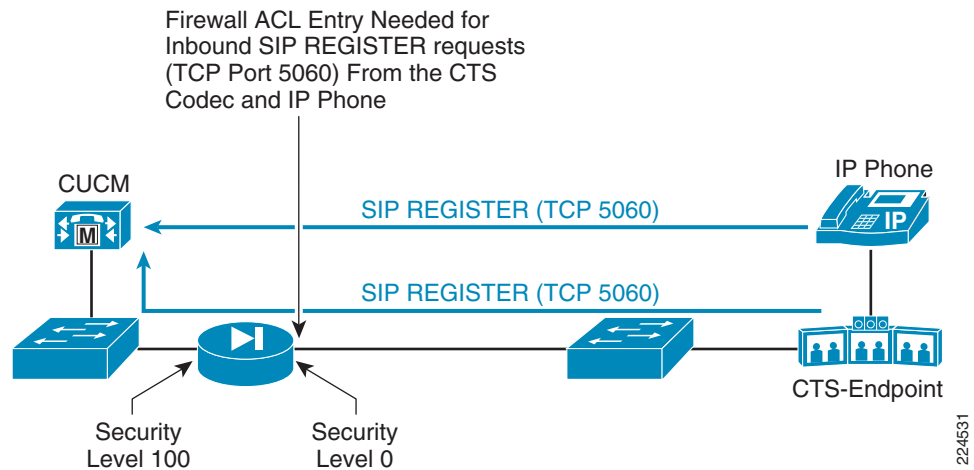
Figure 13-8 shows both protocols required by the CTS endpoints for system image upgrade and configuration download.

Figure 13-8 Configuration Download Protocol Support for TelePresence in a Firewall Deployment

SIP Registration

Once the CTS codec(s) and associated IP phone have completed downloading their configuration files, and possibly upgrading their system images, both perform a SIP registration with the CUCM cluster. SIP signaling uses either TCP or UDP port 5060. The connection-oriented nature of TCP makes it preferred for TelePresence deployments and is the only protocol discussed within this chapter. Note that this section of the discussion also assumes that the CUCM cluster is on the opposite side of the firewall from the CTS endpoint.

Since SIP REGISTER requests are initiated by the CTS primary codec and associated IP phone, ACL entries on the firewall which allow the inbound traffic may be required. Figure 13-9 shows the SIP protocol required for registration by both the CTS primary codec and associated IP phone.

Figure 13-9 SIP Registration Support for TelePresence in a Firewall Deployment

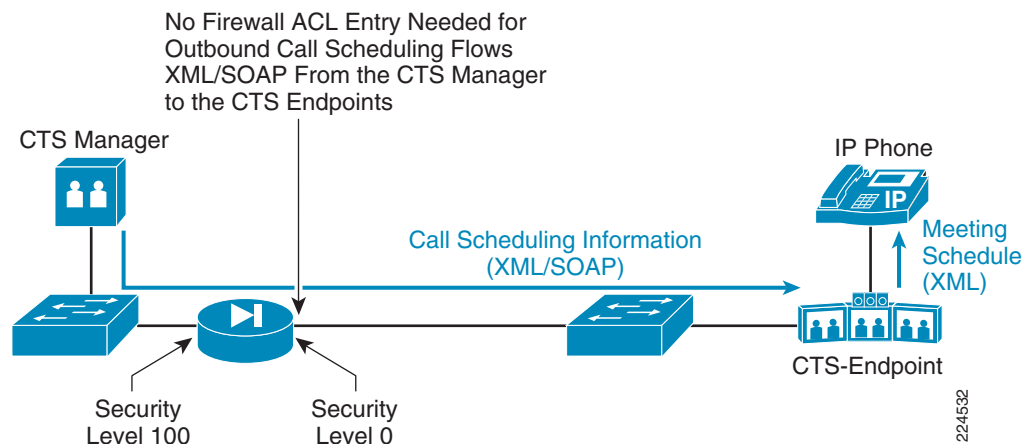
It should be noted that if additional IP phones which support the SCCP protocol exist on the firewall interface with the lower security level, then an ingress ACL entry on the firewall interface with the lower security level may be needed to allow such devices to register with CUCM. These may be utilized due to the audio add-on feature of TelePresence meetings.

Call Scheduling and Services Flows

Call scheduling flows are the data flows between the CTS Manager and CTS endpoints used to update the meeting schedule information which appears on the IP 7970 phone associated with the CTS endpoint. CTS Manager uses XML/SOAP in order to push scheduling information out to the CTS endpoints. CTS endpoints then push the content to the GUI interface of the phone via XML.

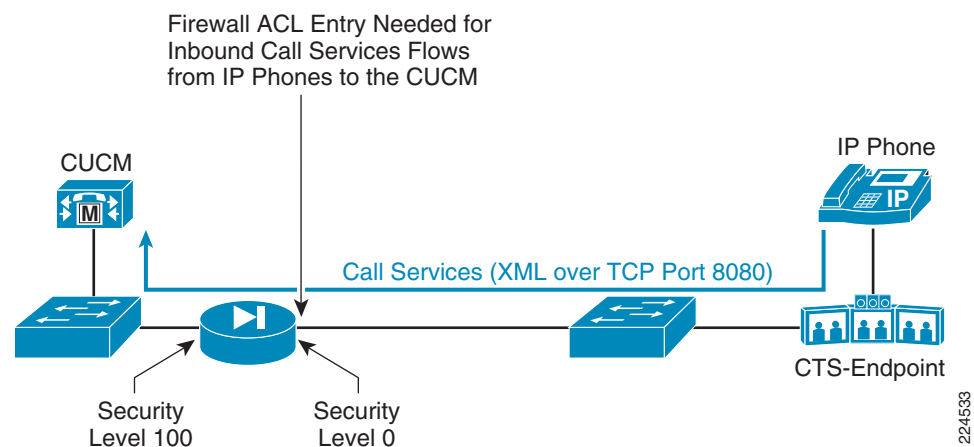
From a firewall perspective, since the session is initiated from the CTS Manager on the higher security level interface to the CTS endpoint on the lower security level interface, no ACL entry is needed to support the call scheduling information flows. This is shown in [Figure 13-10](#).

Figure 13-10 Call Scheduling Flow Support for TelePresence in a Firewall Deployment



Call service flows include things such as directory lookup services available via the graphical user interface of the IP phone associated with the CTS endpoint. IP phones use XML over HTTP port 8080 to communicate with the CUCM cluster and potentially other servers to provide these services. This is shown in [Figure 13-11](#).

Figure 13-11 Call Services Flows Support for TelePresence in a Firewall Deployment

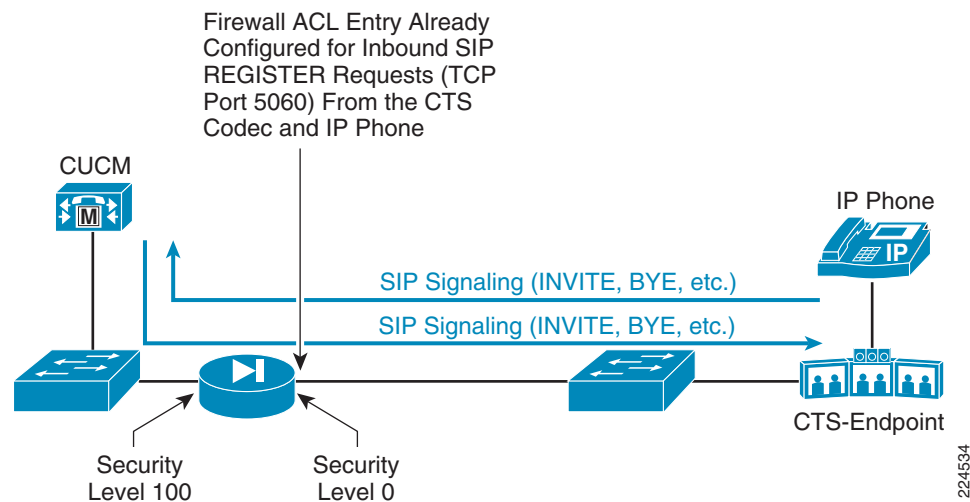


Since these flows originate with the IP phone associated with the CTS endpoint or other IP phones which may be bridged onto a TelePresence meeting, an ingress ACL entry on the firewall interface with the lower security level may be required to allow these to pass through the firewall.

Call Signaling Flows

Call signaling flows are the SIP signaling (UDP and TCP Port 5060) flows between the CTS endpoints and the CUCM used to start, stop, and put calls on hold. From a firewall perspective, since SIP is already allowed through the firewall in order for the CTS endpoint to register with the CUCM cluster, no additional ACL entries are needed to support call signaling flows. This is shown in [Figure 13-12](#).

Figure 13-12 Call Signaling Flow Support for TelePresence in a Firewall Deployment



Note that in [Figure 13-12](#), outbound call signaling from the CUCM to the CTS endpoint on the opposite side of the firewall are automatically allowed from an interface with a higher security level to an interface with a lower security level.

Media Flows

Media flows are the actual RTP/UDP streams that carry the audio and video of the TelePresence meeting. Video streams from individual cameras are carried within individual RTP streams. Likewise, audio streams from individual microphones are carried within individual RTP streams. All of the audio RTP streams are multiplexed into a single audio UDP stream before being sent over the network. Likewise, all of the video RTP streams are multiplexed into a single video UDP stream before being sent over the network. In addition, RTCP control information for each audio and video stream is also multiplexed within each UDP stream. This eases firewall traversal, since only a single UDP audio stream and single UDP video stream is sent from a CTS endpoint.

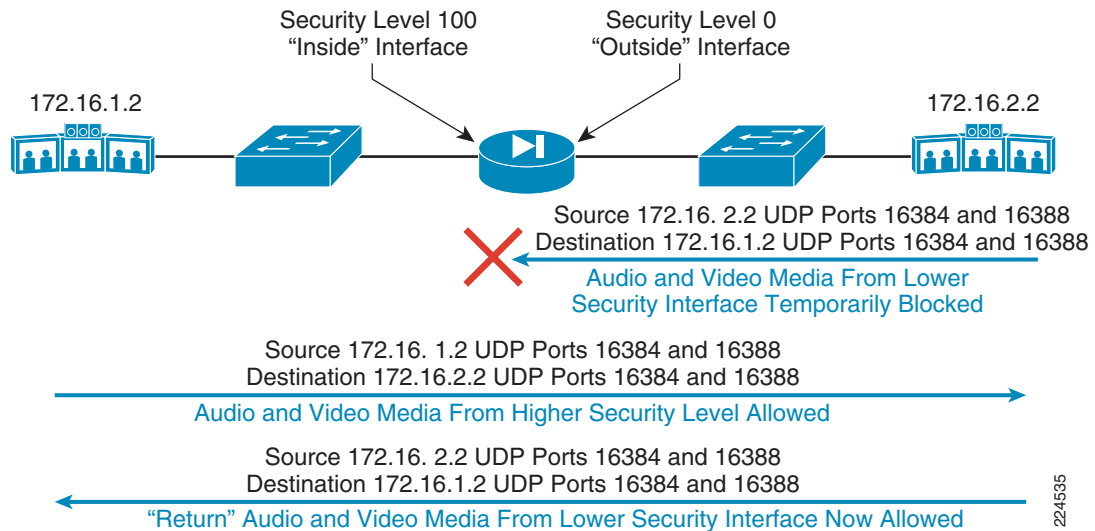
There are slight variations in media flows between point-to-point TelePresence calls and multipoint TelePresence calls. Each is discussed individually.

Point-to-Point TelePresence Calls

In point-to-point TelePresence calls, the audio and video UDP streams flow directly between the CTS endpoints. Because there is only one UDP audio stream and one UDP video stream, the flows are symmetric, as shown in [Figure 13-3](#). Therefore, TelePresence may be deployed without any ACL entry on the lower security interface which allows the inbound UDP media streams (typically UDP ports 16384 and 16388). This assumes no application-level inspection of SIP traffic to dynamically open media ports either. Once SIP signaling has completed, each CTS endpoint independently begins to send

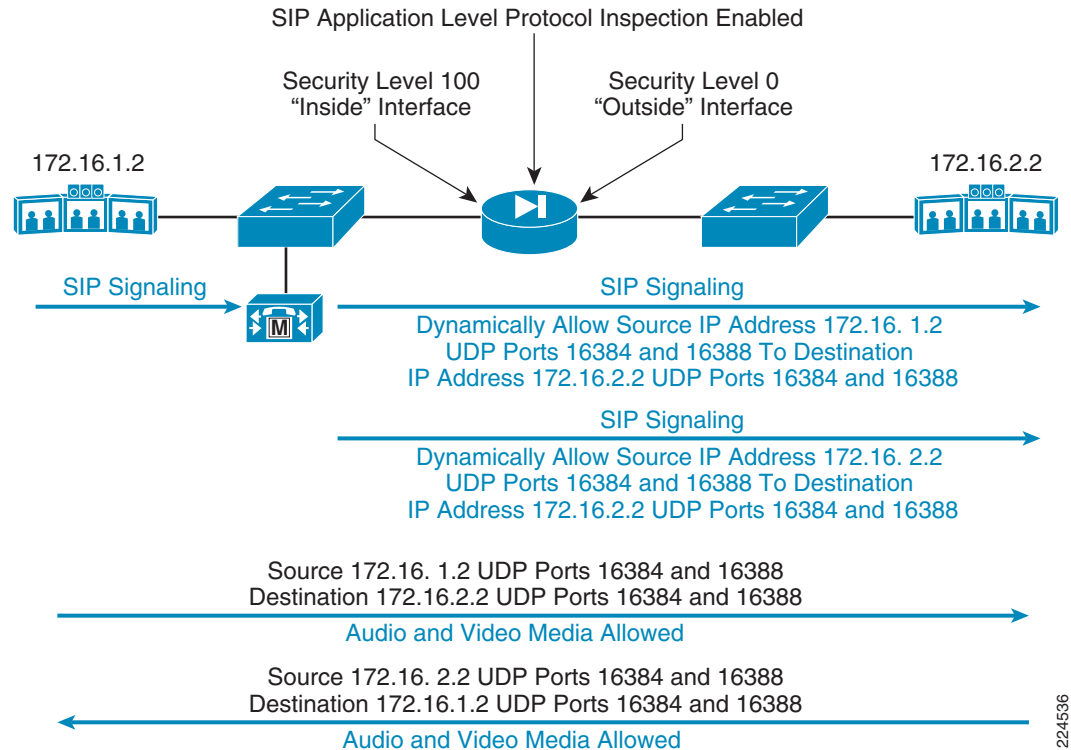
audio and video media. The audio and video media from the CTS endpoint on the lower security interface of the firewall are temporarily blocked by the firewall. However, when the audio and video media from the CTS endpoint on the higher security interface of the firewall passes through the firewall, it will dynamically allow the “reverse” traffic through. This unblocks the audio and video media from the CTS endpoint on the lower security interface of the firewall. Figure 13-13 shows this functionality.

Figure 13-13 Firewall Behavior with Symmetric TelePresence Flows



This behavior allows the network administrator to implement the firewall without having to open static RTP port ranges on the lower security interface, reducing the security vulnerability of the network. The network administrator can certainly configure an ingress ACL on the lower security interface allowing UDP ports 16384 and 16388 (or a larger range of ports from 16384 to 32677) if the initial blocking of video presents an issue.

However, despite this functionality, this configuration is not recommended. Instead, Cisco recommends enabling application-level protocol inspection of SIP traffic in order to allow the firewall to dynamically open and close the necessary UDP ports for the audio and video media. An example of this is shown in Figure 13-14.

Figure 13-14 TelePresence with SIP Application Level Protocol Inspection

Enabling SIP application level protocol inspection causes the firewall to inspect the SDP packets within SIP INVITE requests for the particular IP addresses and UDP ports required for audio and video media in each direction. The firewall continues to inspect any re-INVITES and inspects the SIP BYE message and dynamically closes the media ports as well.

**Note**

Dynamically closing the media ports has been observed to cause temporary blocking of a large number of TelePresence audio and video media packets, since the codecs typically send SIP BYE messages before actually stopping the media flows.

Multipoint TelePresence Calls

In multipoint TelePresence calls, the audio and video UDP streams flow between the CTS endpoints and the CTMS. Again, there is only one UDP audio stream and one UDP video stream for each CTS endpoint. However, since the CTMS has a single IP address and has to support multiple UDP audio and video streams from multiple CTS endpoints, the flows are not necessarily symmetric from a UDP port numbering perspective. Therefore the network administrator may need to statically open a range of UDP ports to the IP address of the CTMS on the lower security interface of the firewall if SIP application-level protocol inspection is not utilized. However, as with point-to-point TelePresence, it is recommended to enable SIP application level protocol inspection in order to allow the firewall to dynamically open and close the necessary media ports.

Note that TelePresence has currently been tested only in a point-to-point configuration within the ESE test lab, both with and without SIP application level protocol inspection enabled.

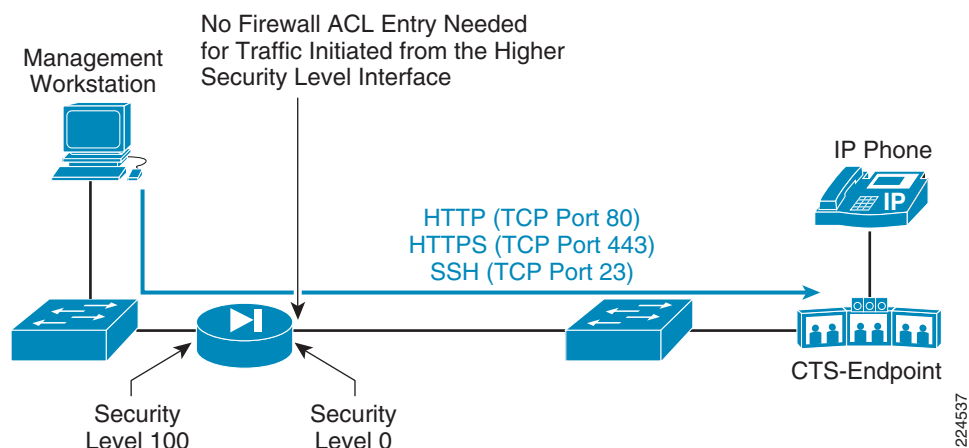
Management Flows

Management flows are needed by management stations or end-user PCs to monitor, configure, and troubleshoot the CTS endpoints. CTS codecs currently support management connections utilizing the HTTPS and SSH protocols. The associated IP phone also supports management connections utilizing HTTP, HTTPS, and SSH protocols. In addition CTS codecs can be configured to send SNMP traps to a management station.

HTTP, HTTPS, and SSH

HTTP utilizes TCP port 80. HTTPS utilizes TCP port 443. SSH utilizes TCP port 23. Since all of these protocols are initiated from management stations on the interface with the higher security level toward the CTS endpoint on the interface with the lower security level, no ingress ACL entry is required on the firewall interface with the lower security level. This is shown in [Figure 13-15](#).

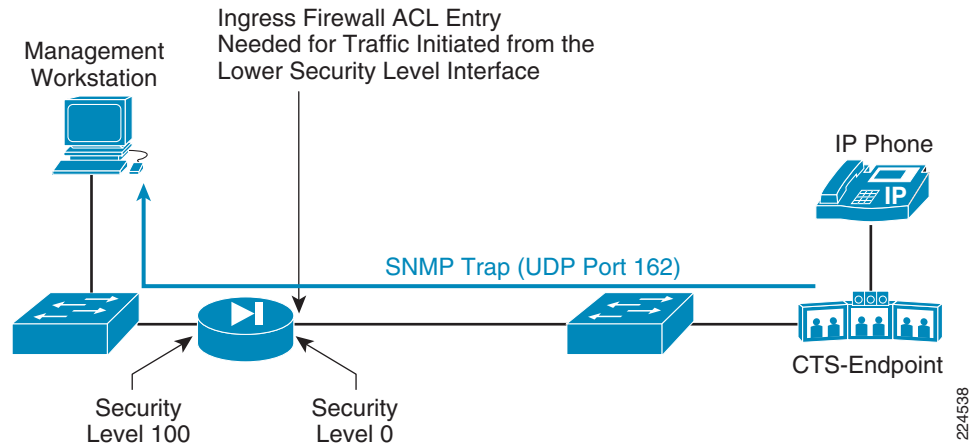
Figure 13-15 Management Station Initiated Protocols



SNMP Traps

SNMP management stations typically listen on UDP Port 162 for SNMP traps. Since SNMP traps are unsolicited events from the CTS endpoints, an ingress ACL entry on the firewall interface with the lower security level may be needed to allow them to pass to the management server on the interface with the higher security level. This is shown in [Figure 13-16](#).

Figure 13-16 CTS Endpoint Initiated Management Protocols



CTS endpoints and their associated IP7970 phones have also been observed to send ICMP Type 3 Code 3 (destination unreachable, port unreachable) packets to each other occasionally. It is not certain why these are sent, perhaps to inform the other side when the audio and video ports have been closed and to stop sending RTP traffic. Operation of TelePresence appears normal with these blocked. However, since they have been observed in firewall syslog output during testing, ICMP Type 3 (destination unreachable) packets may also need to be enabled through the firewall. Enabling SNMP through a firewall should be used with caution however.

ESE Firewall Test Results

Table 13-1 summarizes the results of ESE TelePresence testing in a point-to-point configuration with an ASA 5500 firewall.

Table 13-1 Firewall Configurations Evaluated

Firewall Mode	Interface Security Level	SIP Application Layer Protocol Inspection	One-Sided Address Translation	Results
Routed (Layer 3)	Unequal	Enabled	No NAT	Tested - Passed
			NAT	Tested - Passed
		Disabled	No NAT	Tested - Passed
			NAT	Tested - Failed
	Equal	Enabled	No NAT	Tested - Passed
			NAT	Not Tested
		Disabled	No NAT	Tested - Passed
			NAT	Not Tested

The operation of TelePresence through a firewall operating in transparent mode was not tested. For each of the routed mode firewall configurations tested above, various the call signaling and media flow scenarios were tested, including CTS endpoint reload and upgrade, call initiation from the CTS endpoint on either side of the firewall, and audio add-on of an IP phone on either side of the firewall.

Although for each firewall configuration, both voice and TelePresence traffic were configured into a priority queue on the firewall, no attempt was made at evaluating the performance of the firewall under traffic load. The results of such testing are highly dependent upon the amount and type of traffic traversing the firewall, the number of connections that need to be set up and torn down by the firewall, as well as the structure of the access-control lists allowing or denying traffic. Therefore, only the firewall signaling control plane was evaluated.

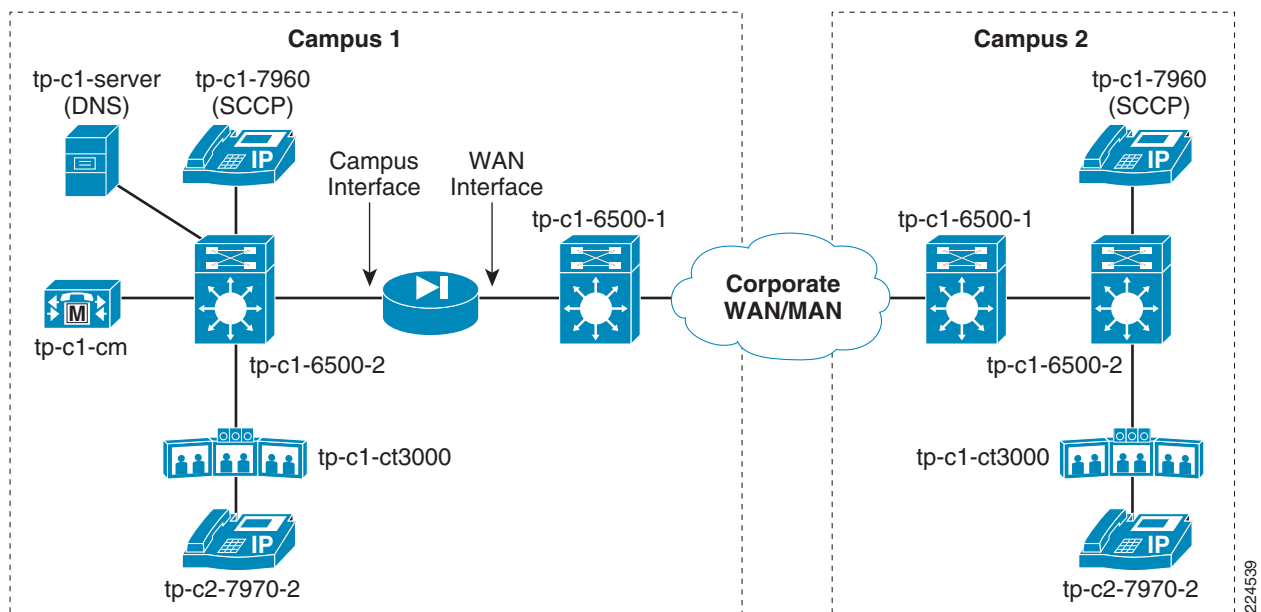
In all the tests with SIP application layer protocol level inspection enabled, only one minor issue was observed. If a TelePresence call is initiated by the CTS endpoint on the opposite side of the firewall from the CUCM, and the call is subsequently put on hold by the CTS endpoint on the same side of the firewall as the CUCM, the SIP signaling between the CUCM and the CTS endpoint on the opposite side of the firewall does not allow the firewall to correctly set up the necessary RTP pinhole connections to allow music-on-hold (MOH) between the CUCM server and the CTS endpoint. Firewall syslog output indicated RTP traffic from the CUCM destined for UDP port 16384 of the CTS endpoint was blocked. No music-on-hold was heard on the CTS endpoint. This is considered to be a minor issue and was only observed because an ingress ACL on the interface with the higher security level was used to determine specific ports required for TelePresence deployments through the firewall. The default behavior of a firewall which allows all flows from an interface with a higher security level to an interface with a lower security level would normally allow this to work correctly.

The test with unequal security levels on the interfaces, one-sided NAT, and SIP application level protocol inspection disabled failed during the CTS endpoint reload test and is therefore not a recommended configuration to support TelePresence.

Example Firewall Configuration

The following shows part of an example firewall configuration between CTS endpoints on two campuses as shown in [Figure 13-17](#).

Figure 13-17 Example Firewall Configuration



The firewall for this example has been configured for routed mode operation, with unequal interface security levels, no NAT, and with SIP application layer protocol inspection enabled.

Interface Security Levels:

```
!
interface GigabitEthernet0/0
  description Connection to tp-c1-6500-2 Gig3/1
  speed 1000
  duplex full
  nameif CAMPUS
  security-level 50! Higher security level
  ip address 10.16.7.2 255.255.255.0
!
!
interface GigabitEthernet1/0
  description Connection to tp-c1-7200-1 Gig0/2
  speed 1000
  duplex full
  nameif WAN
  security-level 40! Lower security level
  ip address 10.16.8.2 255.255.255.0
!
```

Static Xlates:

```
!
static (CAMPUS,WAN) tp-c1-server tp-c1-server netmask 255.255.255.255 dns
! Allow the Campus 1 DNS server to be reachable from the lower security interface
static (CAMPUS,WAN) tp-c1-ct3000 tp-c1-ct3000 netmask 255.255.255.255 dns
! Allow the Campus 1 TelePresence unit to be reachable from the lower security interface
static (CAMPUS,WAN) tp-c1-7960-1 tp-c1-7960-1 netmask 255.255.255.255 dns
! Allow the Campus 1 IP7960 phone to be reachable from the lower security interface
static (CAMPUS,WAN) tp-c1-7970-1 tp-c1-7970-1 netmask 255.255.255.255 dns
! Allow the Campus 1 IP7970 phone to be reachable from the lower security interface
static (CAMPUS,WAN) tp-c1-cm tp-c1-cm netmask 255.255.255.255 dns
! Allow CUCM to be reachable from the lower security interface
```

WAN ACL:

```
!
access-list NEW_WAN extended permit icmp any any unreachable
! Allow ICMP unreachable from all devices to be sent. Not necessarily needed for
TelePresence operation, but
! ICMP unreachable were observed in the firewall syslog output and added to the ACL.
access-list NEW_WAN extended permit tcp host tp-c2-ct3000 host tp-c1-cm eq sip
access-list NEW_WAN extended permit tcp host tp-c2-7970-1 host tp-c1-cm eq sip
access-list NEW_WAN extended permit tcp host tp-c2-7970-2 host tp-c1-cm eq sip
! Allow SIP devices to register with CUCM
access-list NEW_WAN extended permit tcp host tp-c2-7960-1 host tp-c1-cm eq 2000
! Allow SCCP devices to register with CUCM
access-list NEW_WAN extended permit udp host tp-c2-ct3000 host tp-c1-server eq domain
access-list NEW_WAN extended permit udp host tp-c1-7960-1 host tp-c1-server eq domain
access-list NEW_WAN extended permit udp host tp-c1-7970-1 host tp-c1-server eq domain
access-list NEW_WAN extended permit udp host tp-c2-7970-2 host tp-c1-server eq domain
! Allow devices to access the DNS server to translate names to valid IP addresses
access-list NEW_WAN extended permit udp host tp-c2-7970-1 host tp-c1-cm eq tftp
access-list NEW_WAN extended permit udp host tp-c2-ct3000 host tp-c1-cm eq tftp
access-list NEW_WAN extended permit udp host tp-c2-7960-1 host tp-c1-cm eq tftp
access-list NEW_WAN extended permit udp host tp-c2-7970-2 host tp-c1-cm eq tftp
! Allow devices to access the TFTP server within CUCM for downloading of configuration and
OS
access-list NEW_WAN extended permit tcp host tp-c2-7970-1 host tp-c1-cm eq 8080
access-list NEW_WAN extended permit tcp host tp-c2-7960-1 host tp-c1-cm eq 8080
! Allow XML access from the IP Phones to CUCM
```

```

access-list NEW_WAN extended permit tcp host tp-c2-ct3000 host tp-c1-cm eq 6970
! TCP port used during firmware upgrades of the TelePresence CTS-3000 units.
!
! Note: ACL entry for SNMP traps not included in this configuration.
!

```

Campus ACL:

```

!
access-list NEW_CAMPUS extended permit icmp any any unreachable
! Allow ICMP unreachable from all devices to be sent. Not necessarily needed for
TelePresence operation, but
! ICMP unreachable were observed in the firewall syslog output and added to the ACL.
access-list NEW_CAMPUS extended permit udp host tp-c1-cm host tp-c2-ct3000 eq 16384
! Explicitly allows music-on-hold from CUCM to the CTS-3000. Necessary currently because
of a
! potential issue with SIP signaling on the TelePresence CTS-3000s when a call is placed
on hold.
access-list NEW_CAMPUS extended permit tcp host tp-c1-blaster host tp-c2-ct3000 eq https
access-list NEW_CAMPUS extended permit tcp host tp-c1-blaster host tp-c2-ct3000 eq ssh
access-list NEW_CAMPUS extended permit tcp host tp-c1-blaster host tp-c2-7960-1 eq www
access-list NEW_CAMPUS extended permit tcp host tp-c1-blaster host tp-c2-7970-1 eq www
access-list NEW_CAMPUS extended permit tcp host tp-c1-server host tp-c2-ct3000 eq https
access-list NEW_CAMPUS extended permit tcp host tp-c1-server host tp-c2-ct3000 eq ssh
access-list NEW_CAMPUS extended permit tcp host tp-c1-server host tp-c2-7960-1 eq www
access-list NEW_CAMPUS extended permit tcp host tp-c1-server host tp-c2-7970-1 eq www
! This section of allows background and management traffic.
!
! Note: Typically no ACL would be utilized from the higher level security interface to
the lower level security interface.
! The ACL was configured mostly for test purposes. Actual deployments will not likely
utilize ingress ACLs on the
! higher level security interface.

```

Application of ACLs Inbound on Interfaces:

```

!
access-group NEW_CAMPUS in interface CAMPUS
access-group NEW_WAN in interface WAN
!

```

Global Policy Which Includes SIP Application Specific Protocol Inspection:

```

!
policy-map type inspect dns migrated_dns_map_1
  parameters
    message-length maximum 512
!
policy-map asa_global_fw_policy
  class inspection_default
    inspect dns migrated_dns_map_1
    inspect ftp
    inspect h323 h225
    inspect h323 ras
    inspect netbios
    inspect rsh
    inspect rtsp
    inspect http
    inspect esmtp
    inspect sqlnet
    inspect sunrpc
    inspect tftp
    inspect xdmcp
    inspect sip
    inspect skinny

```



```
! SIP, SCCP, TFTP and DNS application layer protocol inspection enabled.  
!  
service-policy asa_global_fw_policy global  
!
```

■ Example Firewall Configuration

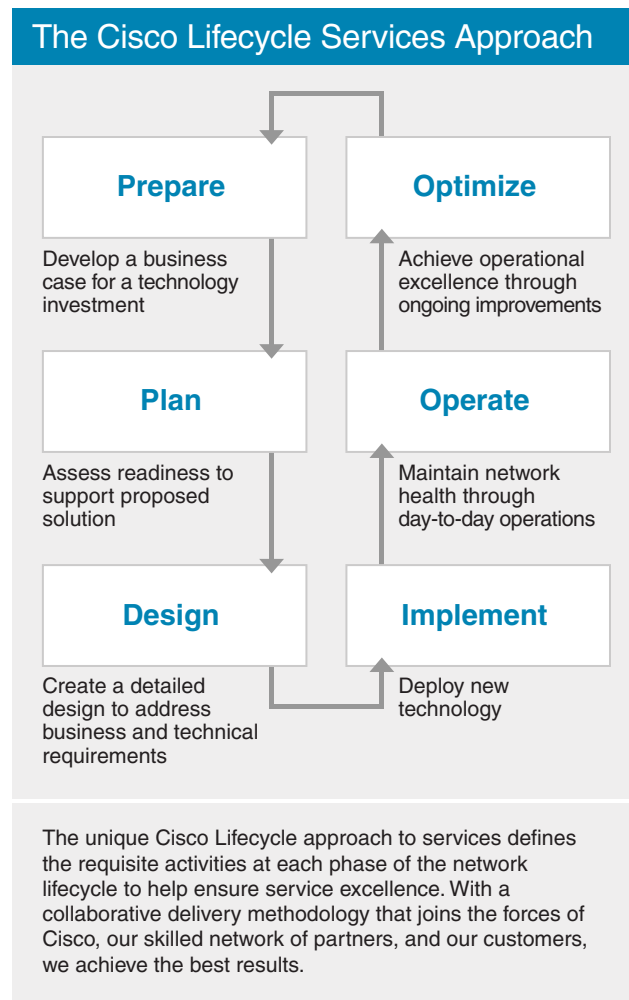


CHAPTER 14

Cisco Services for Cisco TelePresence

Cisco® and its partners provide comprehensive services throughout the planning, design, implementation, and ongoing operation of Cisco TelePresence, helping organizations realize the full potential of the solution.

Figure 14-1 The Cisco Lifecycle Services Approach



Challenge

To compete in today's global business environment, enterprises must be able to effectively communicate, collaborate, and respond rapidly to change—across all geographic boundaries. Cisco TelePresence offers a new technology platform that creates “in-person” experiences between people, places, and events over the IP network. Employees can connect easily and instantly with coworkers, customers, and partners anywhere in the world without leaving the office—speeding decision making, improving business continuity in the event of disasters or disruptions, and providing a distinct competitive edge. However, to gain the full advantages of this solution, organizations need to make sure that Cisco TelePresence is properly deployed and that the critical elements of the solution are functioning optimally at all times.

Solution

Cisco Services for TelePresence provide comprehensive service offerings to help enterprises prepare, plan, and design their networks for the successful implementation of Cisco TelePresence and maintain essential ongoing maintenance and support. These services play an essential role in the successful deployment and ongoing operation of Cisco TelePresence technology by protecting enterprises' Cisco TelePresence investment and helping them achieve the full benefits of the solution. Ultimately, Cisco Services for TelePresence let enterprises focus on business transformation—not the technology.

Cisco Services for TelePresence encompass four service offerings:

- Cisco TelePresence Planning, Design, and Implementation Service
- Cisco TelePresence Essential Operate Service
- Cisco TelePresence Select Operate Service
- Cisco TelePresence Remote Assistance Service

Together, these offerings provide a comprehensive suite of services designed specifically for Cisco TelePresence solutions, based on the Cisco Lifecycle Services framework. The services are available from Cisco or through a set of Cisco Advanced Technology Provider (ATP) partners with deep experience in networking and Cisco Unified Communications and special training in virtual presence technology. These partners draw on proven methodologies to accelerate the business benefits of Cisco TelePresence technology, and focus on the solution 24x7 so that enterprise IT departments can focus on the business.

The Cisco TelePresence Planning, Design, and Implementation Service

Cisco TelePresence technology has a profound effect on an organization's ability to communicate, cooperate, and respond to unforeseen business issues. However, to provide the consistent, high-quality experience enterprise users expect, the organization's network, physical meeting sites, and the Cisco TelePresence solution itself must be optimally designed and implemented. Without careful consideration of an enterprise's specific business and technical requirements, IT and end-user experience, and the effects of Cisco TelePresence on the overall network, organizations might not realize the full potential of the solution.

The Cisco TelePresence Planning, Design, and Implementation Service provides comprehensive support throughout the planning and deployment of a Cisco TelePresence solution, helping organizations quickly realize the benefits of this new real-time, immersive technology. The service helps enterprises achieve their business objectives by assessing the existing network and physical environments, developing an implementation-ready design based on the organization's unique requirements, and working with internal IT staff throughout the implementation and testing of the solution as well as through end-user

training. The service is delivered by expert Cisco or Cisco ATP partner engineers with deep backgrounds in Cisco Unified Communications and a detailed understanding of all components of the Cisco TelePresence solution, including hardware, software, and application configuration. Cisco and its ATP partners draw on the proven Cisco Lifecycle Services methodology, as well as industry-leading Cisco intellectual property and networking expertise to align Cisco TelePresence service and support activities with the enterprise's business and technology requirements throughout the network lifecycle. As a result, enterprises can deploy Cisco TelePresence on their existing network (instead of building an overlay network), help ensure smooth integration and interoperability with other Cisco Unified Communications solutions, and gain maximum advantage from their converged infrastructure investment.

The Cisco TelePresence Planning, Design, and Implementation Service consists of the following service components:

- Cisco TelePresence Prequalification
- Cisco TelePresence Project Management
- Cisco TelePresence Requirements Validation
- Cisco TelePresence Site Survey
- Cisco TelePresence Path Qualification
- Cisco TelePresence Detailed Design Development
- Cisco TelePresence Implementation Plan
- Cisco TelePresence System Acceptance Testing
- Cisco TelePresence End-User Training

In the initial prepare phase of a Cisco TelePresence deployment, Cisco or Cisco partner engineers use the Cisco TelePresence Prequalification Checklist to qualify an enterprise's network and physical meeting locations and verify that they can support the solution. When an enterprise is ready to begin the plan phase of the deployment, Cisco or the Cisco partner delivers a comprehensive Project Management Plan for the implementation and provides a single point of contact for all issues relating to the service. The project team then performs a detailed Requirements Validation to assess the customer's business and technical requirements for the solution and verify that the Cisco TelePresence deployment will meet expectations. This analysis is followed by an exhaustive Site Survey to certify that Cisco TelePresence can operate effectively in the environment. The team then performs Cisco TelePresence Path Qualification—an in-depth examination of the customer network and the links between sites to identify the optimal network path for the solution.

In the design phase of the implementation, the project team develops an implementation-ready Detailed Design for the Cisco TelePresence solution. Then, in the implementation phase, the team develops a comprehensive network implementation plan for each element of the solution, including audiovisual and environmental standards, and deploys the solution. After deployment, the team performs a System Acceptance Testing process, including the creation and implementation of customer-specific test cases for all sites to determine the readiness of the Cisco TelePresence solution for live production. Finally, the project team performs End-User Training, including the development of customized training materials and hands-on education to make sure that system administrators, support staff, and end users all can make full use of Cisco TelePresence technology.

The Cisco TelePresence Essential Operate Service

Even when Cisco TelePresence is expertly deployed, enterprises still need ongoing support and maintenance to safeguard all of the essential components of the solution. However, given the nature of the solution and its many components (hardware and software; voice, video, and data; and room environmental attributes), enterprises need dedicated, system-level support and maintenance to protect their Cisco TelePresence investments.

The Cisco TelePresence Essential Operate Service helps enterprises realize the cost savings and productivity gains that the Cisco TelePresence solution makes possible by delivering a reliable, high-quality meeting experience. Organizations gain 24-hour, 365-day-a-year access to a comprehensive support environment that addresses all aspects of Cisco TelePresence technology—voice and video, software and hardware—with a single, integrated service.

Enterprises gain global, full-time access to highly trained engineers who have a deep understanding of Cisco Unified Communications products and technologies and who specialize in complex IP communications environments. This system-level technical support—whether provided by Cisco or by Cisco certified ATP partners—can help enterprises quickly and cost-effectively resolve issues with any aspect of the Cisco TelePresence solution. If a problem arises with the technology, enterprise IT administrators don't have to determine if the problem lies in the voice, video, or IP aspects of the solution. Instead, this consolidated support model means that one telephone call—or one push of a button on a Cisco Unity® IP Phone—connects administrators with a highly trained technical engineer with deep experience with complex IP communications network issues. These engineers can quickly identify an issue and, if necessary, facilitate collaboration across multiple Unified Communications technology experts to accelerate the resolution of any Cisco TelePresence problem.

Enterprises also gain fast access to replacement parts. The Cisco TelePresence Essential Operate Service includes two Advanced Hardware Replacement options with onsite installation, providing enterprises with parts delivery and replacement by the next business day or within four hours on the same business day, depending on the needs of the organization. The service also includes ongoing operating system and application software updates, strengthening the reliability, functionality, and stability of Cisco TelePresence application software.

In addition, enterprises gain registered access to an array of powerful, industry-leading online support and information systems. These include interactive consulting tools, a comprehensive database, and knowledge transfer resources available through Cisco.com. This robust set of Cisco technical tools and product information increases the self-sufficiency and Unified Communications expertise of internal IT staff, improving productivity while protecting the Cisco TelePresence investment.

The Cisco TelePresence Select Operate Service

While sound planning and ongoing support help enterprises quickly benefit from Cisco TelePresence technology, many organizations do not have the in-house Cisco TelePresence expertise to optimally monitor and manage the solution on a day-to-day basis. Developing that expertise can represent a significant investment in time, people, and resources that can impede the operational efficiency of in-house IT staff.

The Cisco TelePresence Select Operate Service provides 24x7 proactive remote monitoring and management support for the solution to give enterprises greater peace of mind and allow internal IT administrators to focus on core business requirements, instead of Cisco TelePresence. The service combines all of the components of the Cisco TelePresence Essential Operate Service with world-class remote management services from Cisco Remote Operations Services (ROS), including a redundant remote network operations center (NOC) infrastructure that monitors the solution at all times and provides a single point of contact for rapid service restoration. Cisco engineers with in-depth expertise

managing converged infrastructures provide ongoing proactive management, monitoring, reporting, and issue diagnosis and remediation to proactively solve real-time incidents. As a result, enterprises can reduce the operational costs of supporting the solution with in-house resources and help ensure that Cisco TelePresence technology is available, secure, and supporting business goals.

The Cisco TelePresence Remote Assistance Service

Proper planning, expert ongoing support, and real-time remote management can all accelerate the benefits of Cisco TelePresence and give enterprises greater peace of mind using the solution. But what can organizations do to help ensure that virtual presence sessions will always run smoothly? For example, what if the person leading the session schedules an important conference incorrectly or needs to make a change at the last minute and can't remember how?

Enterprises purchasing the Cisco TelePresence Select Operate Service have the option of adding real-time administrative support during Cisco TelePresence sessions with the Cisco TelePresence Remote Assistance Service. By simply pressing the Concierge button on a Cisco Unity IP phone, users in any managed Cisco TelePresence conference room can connect with a Remote Assistance representative day or night, 365 days a year. The service provides remote assistance with scheduling and call setup, answers questions about how to use the solution, and can serve as a single point of contact for any issues requiring engineering support. As a result, users can quickly find answers and resolve unexpected issues, and help ensure a smoother, more effective Cisco TelePresence experience.

Benefits

Cisco Services for TelePresence provide a comprehensive set of activities that are essential to the successful deployment and optimal ongoing operation of Cisco TelePresence technology. The Cisco TelePresence Planning, Design, and Implementation Service protects organizations against downtime caused by improper solution design, helps enterprises avoid costly deployment delays, and helps ensure that the solution fully meets expectations. The Cisco TelePresence Essential Operate Service protects against downtime caused by hardware and software issues and provides critical assistance and expertise to keep innovative Cisco Unified Communications networks running smoothly. Cisco TelePresence Select Operate Service and Remote Assistance Service provide experts that focus on Cisco TelePresence 24x7 so that users and in-house IT administrators can focus on their business. Together, these services deliver a consistent, high-quality Cisco TelePresence experience and allow organizations to focus on transforming their business - not supporting technology.

Cisco Services for TelePresence help organizations:

- Accelerate the business benefits of Cisco TelePresence by accurately assessing the effects of the solution on the network and on physical locations, and addressing potential issues before they arise
- Protect against downtime caused by improper solution design or hardware and software issues
- Decrease deployment times and avoid costly deployment delays, minimizing the risk associated with adopting advanced technologies
- Realize greater peace of mind through proactive remote monitoring and comprehensive operational support and management of all elements of the Cisco TelePresence solution, delivered through a single, dedicated support environment
- Transform business and technical requirements into a detailed design that can be implemented efficiently and effectively and can provide a Cisco TelePresence solution that meets expectations

- Improve the performance and availability of the Cisco TelePresence solution to better meet business requirements and provide a robust foundation for supporting innovative communications applications and intra-company collaboration

Why Cisco Services

Cisco Services make networks, applications, and the people who use them work better together.

Today, the network is a strategic platform in a world that demands better integration between people, information, and ideas. The network works better when services, together with products, create solutions aligned with business needs and opportunities.

The unique Cisco Lifecycle approach to services defines the requisite activities at each phase of the network lifecycle to help ensure service excellence. With a collaborative delivery methodology that joins the forces of Cisco, our skilled network of partners, and our customers, we achieve the best results.

For More Information

For more information about the Cisco Services for TelePresence or other Cisco services, visit www.cisco.com/go/telepresenceservices or contact your Cisco service account manager.