



Cisco CallManager Clusters

This chapter discusses the concept, provisioning, and configuration of Cisco CallManager clusters. Clusters, which were introduced with Cisco CallManager Release 3.0, provide a mechanism for distributing call processing seamlessly across a converged IP network infrastructure to support IP telephony, facilitate redundancy, and provide feature transparency and scalability.

This chapter discusses the operation of clusters within both campus and WAN environments and proposes reference designs for implementation. The following sections cover these topics:

- Cluster Operation and Scalability Guidelines, page 3-1
- Cisco CallManager Redundancy, page 3-6
- Campus Clustering Guidelines, page 3-12
- Intercluster Communication, page 3-14
- Intracluster and Intercluster Feature Transparency, page 3-21

Cluster Operation and Scalability Guidelines

With Cisco CallManager Release 3.0(5), a cluster can contain as many as eight servers, of which six are capable of call processing. The other two servers can be configured as a dedicated database publisher and a dedicated TFTP server, respectively.

The database publisher is used to make all configuration changes and also to produce call detail records. The TFTP server facilitates the downloading of configuration files, device loads (operating code), and ring types.

A dedicated database publisher and a dedicated TFTP server are recommended for large systems. For smaller systems, the function of database publisher and the TFTP server can be combined. Table 3-1 provides guidelines for scaling devices with Cisco CallManager clusters.

Table 3-1 Cisco CallManager Cluster Guidelines

Required Number of IP Phones within a Cluster	Recommended Number of Cisco CallManagers	Maximum Number of IP Phones per Cisco CallManager
2,500	Three servers total: <ul style="list-style-type: none"> • Combined publisher / TFTP • One primary Cisco CallManager • One backup Cisco CallManager 	2,500
5,000	Four servers total: <ul style="list-style-type: none"> • Combined publisher / TFTP • Two primary Cisco CallManagers • One Backup Cisco CallManager 	2,500
10,000	Eight servers total: <ul style="list-style-type: none"> • Database publisher • TFTP server • Four primary Cisco CallManagers • Two backup Cisco CallManagers 	2,500

The preceding recommendations provide an optimum solution. It is possible to reduce the amount of redundancy, and hence use fewer servers. For small systems the database publisher, TFTP server, and Cisco CallManager backup functions can be combined.

The maximum number of registered devices per Cisco CallManager is 5000 in the case of the MCS-7835, including a maximum of 2500 IP telephones, gateways, and Digital Signaling Processor (DSP) devices such as transcoding and conferencing resources. In the event of failure of one of the Cisco CallManagers within the cluster, the maximum number of registered devices remains 5000 per Cisco CallManager in the case of the MCS-7835.

Device Weights

Many types of devices can register with a Cisco CallManager. Each of these resources—IP phones, voice mail ports, Telephony Application Programming Interface (TAPI) devices, Java Telephony API (JTAPI) devices, gateways, and DSP resources such as transcoding and conferencing—carries a different weight. Table 3-2 shows the weight for each of the resource types, based on the consumption of memory and CPU resources.

Table 3-2 Weights by Device Type

Device type	Weight per Session/ Voice Channel	Session/DS0 per Device	Cumulative Device Weight
IP phone	1	1	1
Analog gateway ports	3	Varies	3 per DS0
T1 gateway	3	24	72 per T1
E1 gateway	3	30	90 per E1
Transcoding resource	3	Varies	3 per session
Software MTP	3	48	144 ¹
Conference resource (hardware)	3	Varies	3 per session
Conference resource (software)	3	48	144 ¹
CTI port (TAPI and JTAPI)	20	1	20
Cisco SoftPhone	20	1	20
Messaging (voice mail)	3	Varies	3 per session
Intercluster trunk	3	Varies	3 per session

1. When installed on the same server as Cisco CallManager, the maximum number of sessions is 48.

The total number of device units that a single Cisco CallManager can control depends on the server platform. Table 3-3 gives details of the maximum number of devices per platform.

Table 3-3 Maximum Number of Devices per Server Platform

Server Platform Characteristics	Maximum Device Units per Server	Maximum IP Phones per Server
MCS-7835-1000 ¹ PIII 1000MHz, 1G RAM	5000	2500
MCS-7835 PIII 733MHz, 1G RAM	5000	2500
MCS-7830 PIII 500MHz, 1G RAM	3000	1500
MCS-7830 PIII 500MHz, 512M RAM	1000	500
MCS-7825-800 ¹ PIII 800MHz, 512M RAM	1000	500
MCS-7822 PIII 550MHz, 512M RAM	1000	500
MCS-7820 PIII 500MHz, 512M RAM	1000	500

1. This server platform will not be available until first quarter of 2001.

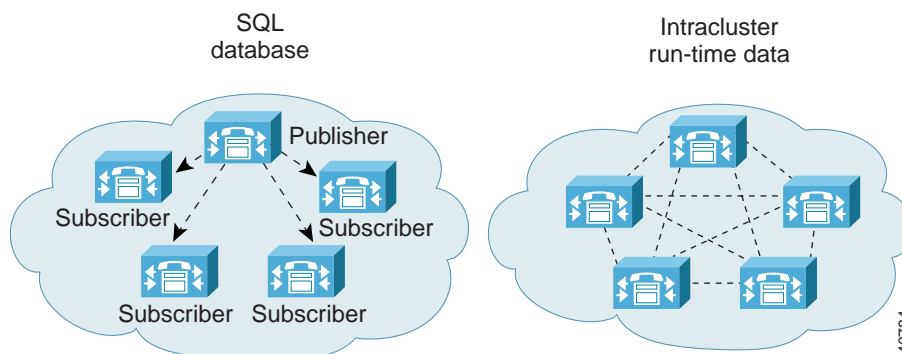
The total number of IP phones that can register with a single Cisco CallManager is limited to 2500 on an MCS-7835, even if only IP phones are registered. To calculate the number of IP phones you can register with a Cisco CallManager, subtract the weighted value of non-IP phone resources from the maximum number of device units allowed for that platform. In the case of the MCS-7835, the maximum number of device units is 5000.

Intracluster Communication

There are two primary kinds of intracluster communications within a Cisco CallManager cluster (Figure 3-1). The first is a mechanism for distributing the database that contains all the device configuration information. The configuration database (Microsoft SQL 7.0) is stored on a publisher and replicated to the subscriber members of the cluster. Changes made on the publisher are communicated to the subscriber databases, ensuring that the configuration is consistent across the members of the cluster as well as facilitating spatial redundancy of the database.

The second intra-cluster communication is the propagation and replication of run-time data such as registration of IP phones, gateways, and DSP resources. This information is shared across all members of a cluster and assures the optimum routing of calls between members of the cluster and associated gateways.

Figure 3-1 Intracluster Communications

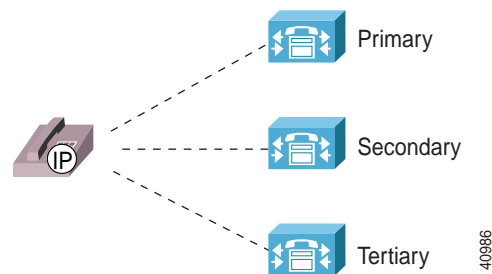


Cisco CallManager Redundancy

Within a cluster, each registered IP phone can be assigned a prioritized list of up to three Cisco CallManagers with which it can register for call processing. Shared resources such as gateways using the Skinny Gateway Protocol are also capable of using this redundancy scheme. The Media Gateway Control Protocol (MGCP) also operates in a similar fashion to provide spatial redundancy for call processing. Peer-to-peer protocols such as H.323 also facilitate redundancy.

Figure 3-2 depicts the redundancy scheme using three Cisco CallManagers.

Figure 3-2 Cisco CallManager Redundancy Group



Each IP phone maintains active TCP sessions with its primary and secondary Cisco CallManagers. This configuration facilitates switchover in the event of failure of the primary Cisco CallManager. Upon restoration of the primary, the device reverts to its primary Cisco CallManager.

Redundancy Group Configurations

You can design your system to provide call processing redundancy by configuring Cisco CallManager redundancy groups. A Cisco CallManager redundancy group is a prioritized list of up to three Cisco CallManagers. You can then assign individual devices to a specific Cisco CallManager redundancy group. A Cisco CallManager redundancy group is a subset of a cluster; all members of a redundancy group are also members of a cluster.

**Note**

The sizes of clusters and redundancy groups are subject to change in future releases of Cisco CallManager.

The following recommendations apply to the configuration of redundancy groups for Cisco CallManager Release 3.0(5):

- Cisco CallManager cluster for up to 2500 users:
 - Server A is a dedicated database publisher and TFTP server.
 - Server B is the primary Cisco CallManager for all registered devices.
 - Server C is the backup Cisco CallManager for all registered devices.

In the configuration above, only a single Cisco CallManager redundancy group is required for servers B and C.

- Cisco CallManager cluster for up to 5000 users:
 - Server A is a dedicated database publisher and TFTP server.
 - Server B is the primary Cisco CallManager for IP phones 1 through 2500.
 - Server C is the primary Cisco CallManager for IP phones 2501 through 5000.
 - Server D is the backup Cisco CallManager for all registered devices.

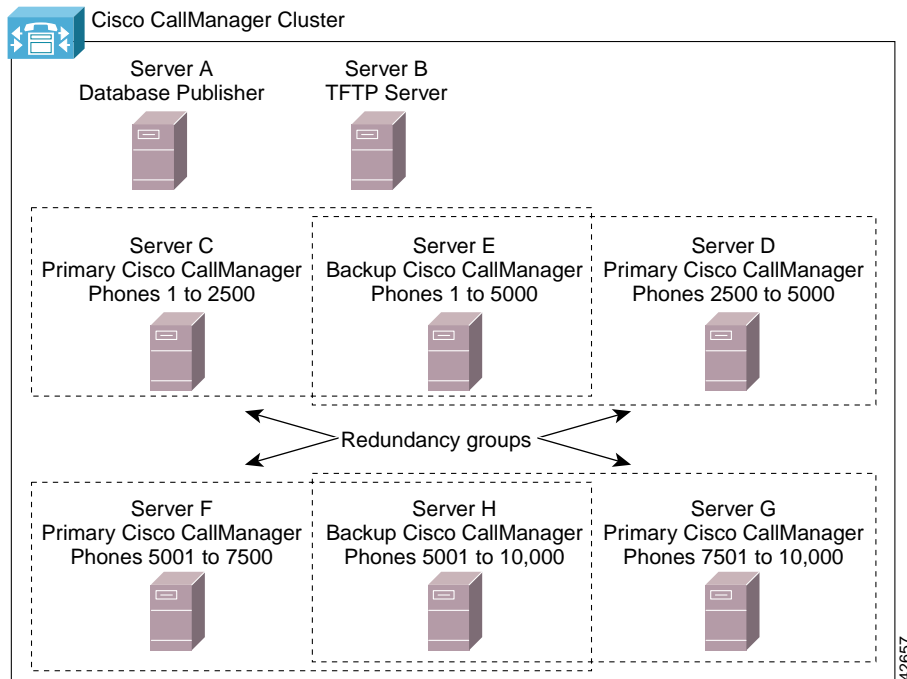
In the configuration above, two Cisco CallManager redundancy groups are required for servers BD and CD.

- Cisco CallManager cluster for up to 10,000 users:
 - Server A is a dedicated database publisher.
 - Server B is a dedicated TFTP server.
 - Server C is the primary Cisco CallManager for IP phones 1 through 2500.
 - Server D is the primary Cisco CallManager for IP phones 2501 through 5000.
 - Server E is the backup Cisco CallManager for IP phones 1 through 5000.
 - Server F is the primary Cisco CallManager for IP phones 5001 through 7500.

- Server G is the primary Cisco CallManager for IP phones 7501 through 10,000.
- Server H is the backup Cisco CallManager for IP phones 5001 through 10,000.

In the above configuration, four Cisco CallManager redundancy groups are required for servers CE, DE, FH and GH. Figure 3-3 illustrates this configuration. Triple redundancy is also possible in this case by configuring the redundancy groups as CEH, DEH, FHE and GHE.

Figure 3-3 Redundancy Groups for a Large System



Note

In the event of a Cisco CallManager failure, calls can be dropped and might need to be reestablished.

Device Pool Configuration

You can use device pools to scale and simplify the distribution of Cisco CallManager redundancy groups. A device pool allows you to assign the following three primary attributes globally to devices:

- Region—Required only if multiple voice codecs are used within an enterprise.
- Date/time group—Specifies date and time zone for a device.
- Cisco CallManager redundancy group—Specifies a list of up to three Cisco CallManagers, which can be used for call processing in a prioritized list.

Figure 3-4 shows an example of a device pool configuration screen. The calling search space for auto-registration is relevant only if auto-registration of IP phones is enabled. This can be used, for example, to limit access of the PSTN to auto-registered devices. Auto-registration is a valuable tool for the initial provisioning of IP phones.

Figure 3-4 Device Pool Configuration Screen

In Figure 3-4, a device pool called Branch 1 G.711 ADE is configured with the following characteristics:

- It is assigned the region Branch 1 G.711. This region contains devices that are capable of communicating by means of G.711 only, such as a voice mail system or conference bridge.
- It is assigned to the appropriate date/time group.
- It is assigned the Cisco CallManager redundancy group ADE, where Cisco CallManager A is the primary, D is the secondary, and E is the tertiary.

A second device pool, called Branch 1 G.729 ADE, could be configured with the following characteristics:

- It is assigned the region Branch 1 G.729. This region contains devices that are capable of communicating by means of both G.729 and G.711, such as IP phones.
- It is assigned to the appropriate date/time group.
- It is assigned the Cisco CallManager redundancy group ADE, where Cisco CallManager A is the primary, D is the secondary, and E is the tertiary.

The same Cisco CallManager group is used for both device pools. However, it is now possible to specify interregion communication codec requirements:

- Intraregion communication uses G.711.
- Interregion communication uses G.729 across the WAN.
- All calls to the G.711 region use G.711. This is required, for example, when accessing an application that is G.711 only.
- This configuration is depicted in Figure 3-5.

Figure 3-5 Interregion Configuration Screen

Region Configuration

Region: Branch 1 G.729
Status: Ready

New Update Delete Cancel

Region Name* Branch 1 G.729

The acceptable rate within this region and between 2 other regions is :

Branch 1 G.729 - Branch 1 G.729	G.711	kbps
Branch 1 G.729 - Branch 1 G.711	G.711	kbps
Branch 1 G.729 - Default	G.729	kbps

* indicates required item

40788

The exact clustering model—and hence device pools used—is driven by the deployment model. The typical device pool configurations, however, have the following characteristics:

- Single-site cluster with no WAN voice interconnectivity
 - Device pools are configured based only on Cisco CallManager redundancy groups.
- Single-site cluster with WAN voice interconnectivity
 - Device pools are configured as above, but with the addition of regions for codec selection. Each cluster could have a G.711 and G.729 region per Cisco CallManager redundancy group.
 - Total device pools = regions x Cisco CallManager redundancy groups.
- Multi-site WAN with centralized call processing
 - Only a single Cisco CallManager redundancy group exists. However, a G.711 region and a G.729 region are required per location. This permits, for example, intrabranch calls to be placed as G.711 and interbranch calls to be placed as G.729.
 - Total device pools = number of sites x regions.

Campus Clustering Guidelines

All members of a Cisco CallManager cluster must be interconnected over a LAN. Cisco CallManager Release 3.0(5) clusters are not supported over a WAN.

The following considerations apply when configuring a campus IP telephony network:

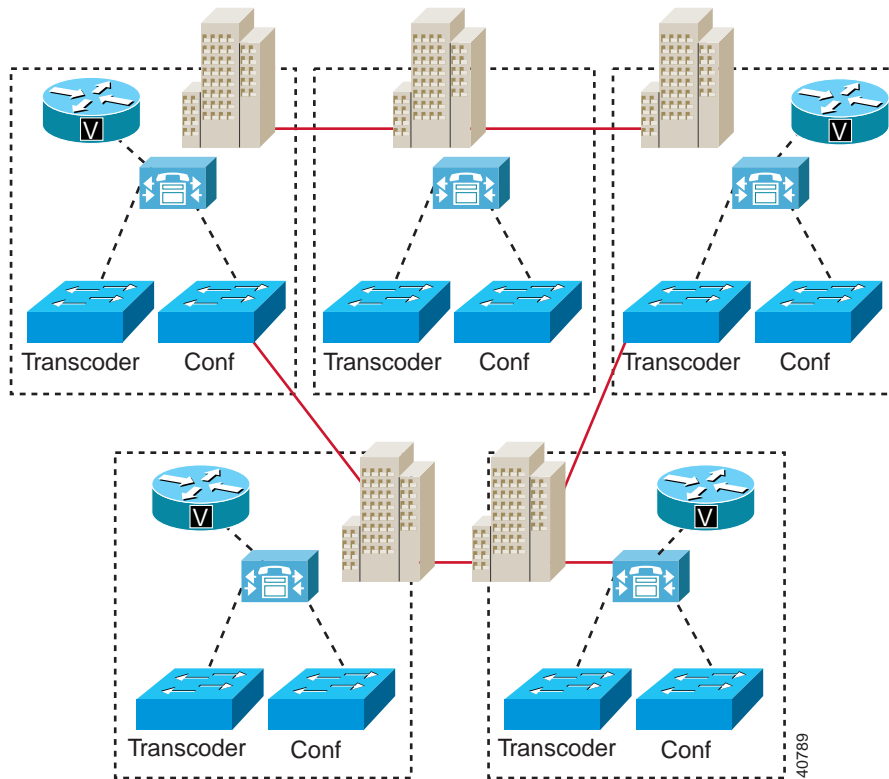
- Maximum of eight servers per cluster with Cisco CallManager release 3.0(5).
- Maximum of 10,000 total registered devices.
- Maximum of 2500 registered IP phones or 3000 devices per Cisco CallManager, including devices registered under failure conditions.
- Switched infrastructure to the desktop (shared media is not supported).

Within a switched campus infrastructure, you can generally assume that the bandwidth is adequate for voice applications. This bandwidth availability depends upon appropriate design and capacity planning within the campus in addition to the establishment of a trust boundary and the required queuing, as discussed in Chapter 2, “Campus Infrastructure Considerations.” There is no requirement for call admission control within a campus cluster.

Cisco CallManager servers should be distributed within the campus to provide spatial redundancy and resiliency. Many metropolitan sites and campus buildings may have only a single conduit providing IP connectivity to other members of the cluster. In this case, if IP connectivity fails, local call processing must be maintained by means of a local server. Gateway resources for PSTN access should likewise be placed strategically to provide the highest possible availability.

Figure 3-6 depicts a typical campus or metropolitan-area network (MAN) cluster deployment.

Figure 3-6 Campus or MAN Cluster



In Figure 3-6 a Cisco CallManager is placed at each of the five buildings or sites. This configuration ensures that, in the event of a failure, local call processing is possible at each site. In cases where diverse routing of fiber cable negates the requirement for a local Cisco CallManager, all call processing could be located in one or more data centers.

Resources such as transcoding and the conferencing DSP are not shared resources and must be provisioned per Cisco CallManager. Once again, where fault tolerance is required, these resources require duplication, and spatial redundancy is recommended. This can be achieved by positioning these resources in strategically placed multi-layer switches.

Intercluster Communication

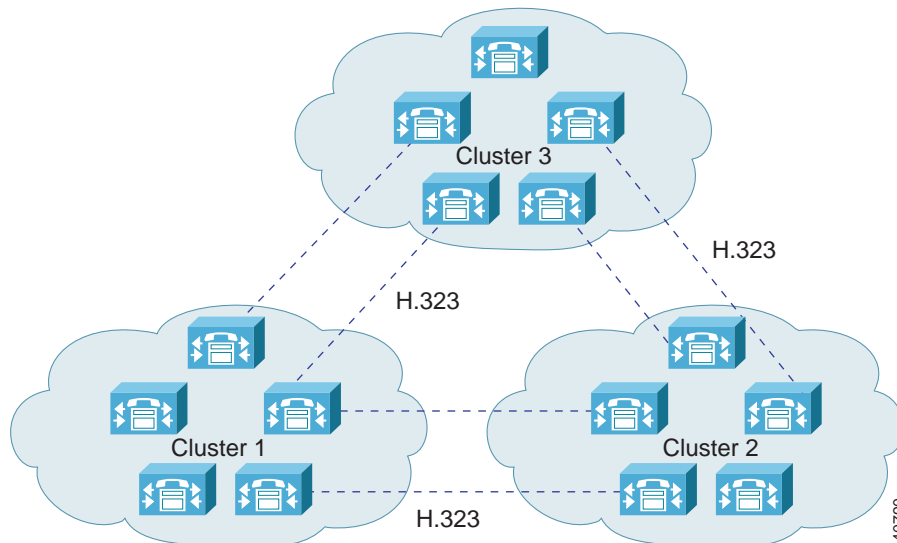
The following sections discuss intercluster communications and address issues in cluster provisioning for isolated campus deployment, multisite WAN deployment with distributed call processing, and multisite WAN deployment with centralized call processing.

Cluster Provisioning for the Campus

Where the requirement for a campus network exceeds 10,000 users, additional clusters are required. Similarly, if local call processing in each site or building requires more than the maximum number of Cisco CallManagers permitted in one cluster, additional clusters are needed.

Communication between clusters is achieved using standards-based H.323 signaling. With a large campus or MAN, where bandwidth is typically over-provisioned and under-subscribed, intercluster call admission control is not required. Figure 3-7 demonstrates this connectivity between clusters within a local area environment.

Figure 3-7 Campus Intercluster Communication Using H.323



In Figure 3-7 the dotted lines represent the H.323 intercluster links, which are configured in pairs to provide redundancy in the event of loss of IP connectivity to any member of the cluster. If desired, you could configure these links as a full mesh. However, Cisco recommends limiting intercluster configuration to two peers. In the majority of situations, this is sufficient to provide adequate resiliency. For deployments where a gatekeeper is used, Cisco recommends a single H.323 connection per cluster. You can implement redundancy by using a Cisco CallManager redundancy group assigned to the gatekeeper.

Unlike earlier releases of Cisco CallManager, release 3.0(5) does not require the use of an MTP to allow supplementary services for H.323 devices. Cisco CallManager 3.0(5) uses the “empty capabilities set” of H.323v2 to facilitate the opening and closing of logical channels between H.323 devices such as Cisco CallManager clusters and Cisco IOS gateways running Cisco IOS Release 12.0(7)T or greater.

Clusters for Multisite WAN with Distributed Call Processing

Where clusters are interconnected over a WAN, there is a pinch point for congestion between clusters, and the network should be engineered to accommodate the required volume of voice traffic. In such cases a method of providing call admission control is required. Because clusters are interconnected using H.323, a Cisco IOS gateway can be added to facilitate this gatekeeper function. Each cluster can be designated as a zone with a maximum configured bandwidth for voice calls.

When using a gatekeeper, Cisco CallManager requests 128 kbps of bandwidth per G.711 inter-cluster call and 20 kbps of bandwidth per G.729a intercluster call. In general, Cisco recommends configuring a single codec for calls that traverse the WAN because this greatly simplifies the provisioning of bandwidth.

Table 3-4 and Table 3-5 give recommendations for bandwidth configuration for intercluster calls.

Table 3-4 Recommended Bandwidth Configuration for Intercluster Calls Using G.729

Number of Intercluster Calls	Bandwidth Required per Call		Bandwidth Required on WAN Links (LLQ/CBWFQ ¹)		Bandwidth Configured on Gatekeeper	
	Without cRTP ²	With cRTP	Without cRTP	With cRTP	Without cRTP	With cRTP
2	24 kbps	12 kbps	48 kbps	24 kbps	40 kbps	40 kbps
5	24 kbps	12 kbps	120 kbps	60 kbps	100 kbps	100 kbps
10	24 kbps	12 kbps	240 kbps	120 kbps	200 kbps	200 kbps

1. Low latency queuing/class based weighted fair queuing
2. Compressed Real-time Transport Protocol

Table 3-5 Recommended Bandwidth Configuration for Intercluster Calls Using G.711

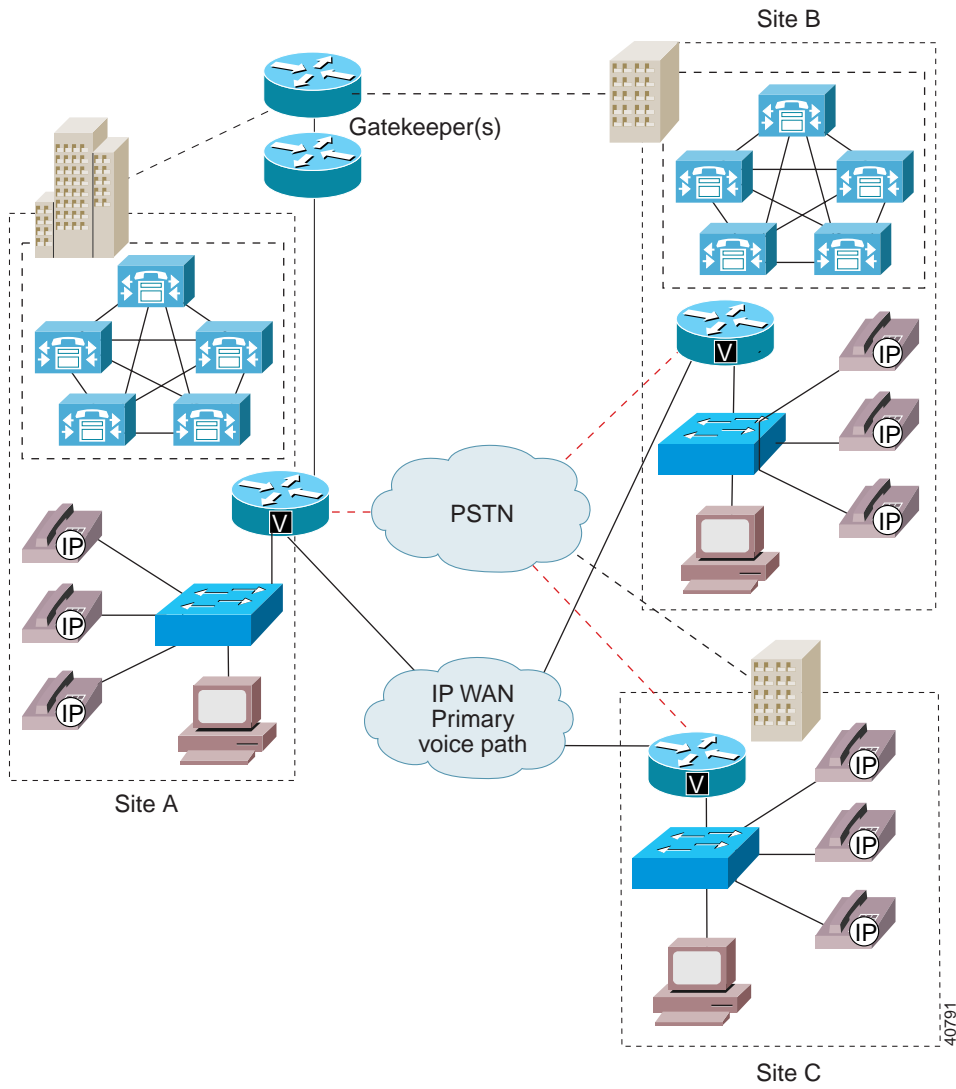
Number of Intercluster Calls	Bandwidth Required per Call	Bandwidth Required on WAN Links (LLQ/CBWFQ)	Bandwidth Configured on Gatekeeper
2	80 kbps	160 kbps	256 kbps
5	80 kbps	400 kbps	640 kbps
10	80 kbps	800 kbps	1280 kbps

The use of gatekeepers provides both inbound and outbound call admission control. With Cisco CallManager Release 3.0(5), a maximum of 100 Cisco CallManagers can register with a gatekeeper. This method of call admission control is restricted to a single active gatekeeper per network. Redundancy can be achieved using the Hot Standby Routing Protocol (HSRP) between two gatekeepers.

Gatekeeper call admission control is a policy-based scheme. It requires static configuration of available resources and is not aware of network topology. It is, therefore, necessary to restrict gatekeeper call admission control schemes to hub-and-spoke topologies with the redundant gatekeeper or gatekeepers (using HSRP) located at the hub. The WAN must be provisioned accordingly, and the voice priority queue must be dimensioned to support all admitted calls.

Figure 3-8 illustrates this deployment model.

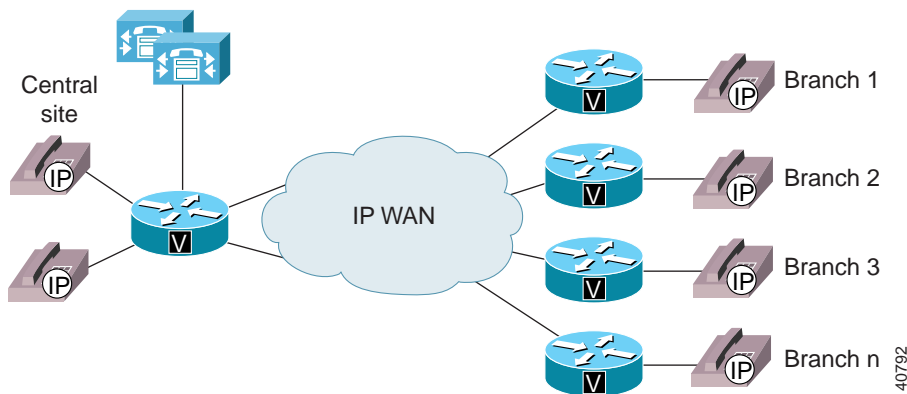
Figure 3-8 Intercluster Communication Using Gatekeepers



Clusters for Multisite WAN with Centralized Call Processing

As stated earlier, Cisco CallManagers within a cluster must be interconnected over a local area network. Cisco CallManager also provides locations-based call admission control that enables provisioning of small branch and telecommuter solutions where remote call processing is acceptable. Figure 3-9 illustrates this model.

Figure 3-9 Locations Based Call Admission Control



In the scheme depicted in Figure 3-9, call processing is maintained only at the central site, and the devices at the branches are configured as belonging to a location. For example, branch 1 might have 12 IP phones, each configured to be in the location Branch 1. Cisco CallManager is then able to track the used and unused bandwidth per location, and admit or deny WAN calls accordingly.

This scheme has been expanded with Cisco CallManager Release 3.0(5) to allow centralized call processing for as many as 2500 remote devices. To implement this type of solution with Cisco CallManager Release 3.0(5), a dedicated Cisco CallManager cluster is required with a single active Cisco CallManager to maintain call state and call admission control.



Note

In this type of centralized configuration, there is a maximum of 2500 users per cluster, regardless of the number of Cisco CallManagers in the cluster (1, 2, or 3 for redundancy purposes). In addition, only one Cisco CallManager in the centralized cluster can be active at a time.

To ensure that only a single Cisco CallManager is active at a time, all devices should be assigned to a single Cisco CallManager redundancy group. This Cisco CallManager redundancy group consists of a prioritized list of up to three Cisco CallManagers. For a centralized call processing cluster, only a single Cisco CallManager redundancy group is recommended, and it should be the default group. In the example shown in Figure 3-10, the redundancy group consists of three Cisco CallManagers, with A as the primary, B the secondary, and C the tertiary Cisco CallManager.

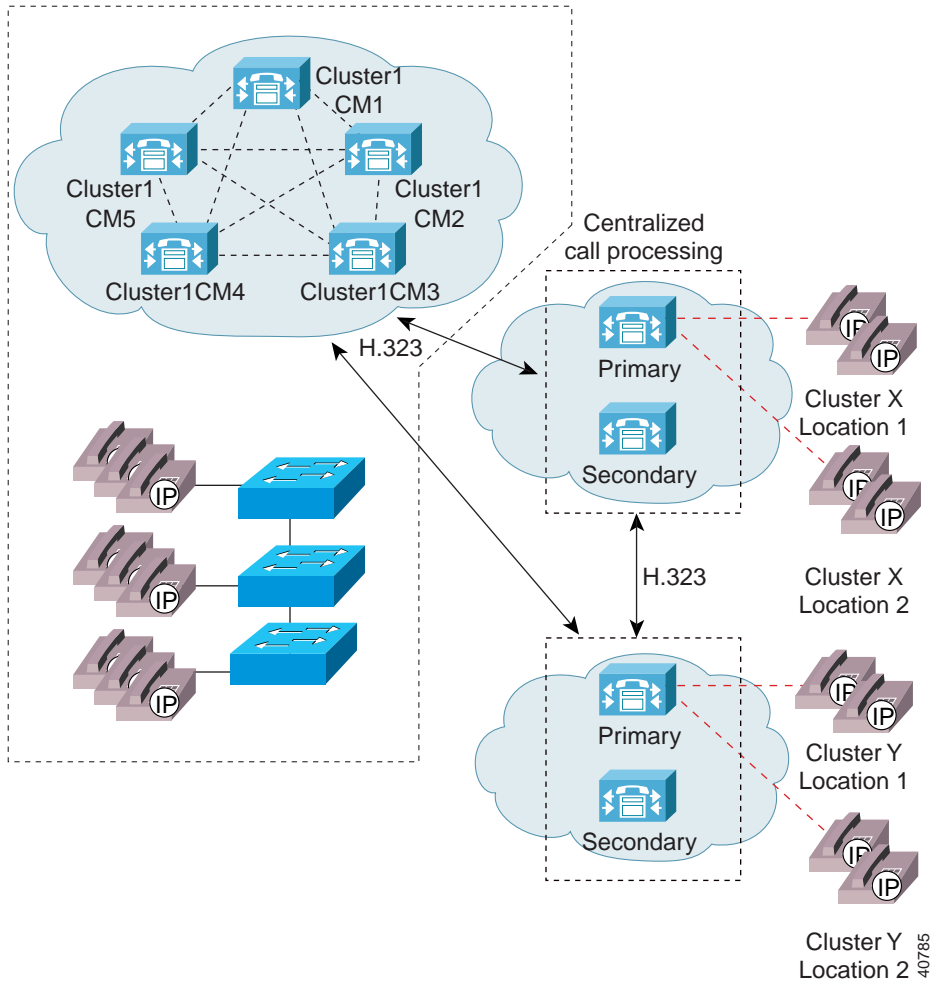
Figure 3-10 Cisco CallManager Redundancy Group Configuration



A typical centralized call processing model might deploy only two Cisco CallManagers. In this case, Cisco recommends that the normally inactive (secondary) Cisco CallManager be the publisher. For a cluster of three Cisco CallManagers, we recommend a dedicated publisher (tertiary) with IP phones and gateways assigned to the primary and secondary Cisco CallManagers.

Figure 3-11 depicts a hybrid deployment model in which a campus cluster is interconnected with two clusters that perform centralized call processing. This example shows that multiple centralized call processing clusters can be deployed and interconnected using H.323. Connectivity to the campus cluster is also achieved using H.323. If intercluster call admission control is required, a gatekeeper can be assigned.

Figure 3-11 Centralized Call Processing Cluster Interconnected with Two Clusters



Intracuster and Intercluster Feature Transparency

The distributed architecture of a Cisco CallManager cluster provides the following primary benefits for call processing:

- Spatial redundancy
- Resiliency
- Availability
- Survivability

In addition, a cluster provides transparent support of user features across all devices in the cluster. This enables distributed IP telephony to span an entire campus or high-speed metropolitan-area network (MAN) with full features.

Intercluster communication provided by H.323 permits a subset of the features to be extended between clusters. These features are currently available between clusters:

- Basic call setup
- G.711 and G.729 calls
- Multiparty conference
- Call hold
- Call transfer
- Calling line ID

In addition, Call Park is available within a cluster but not between clusters.

