

Managing Hardware With Oracle® Solaris Cluster 4.4

ORACLE®

Part No: E69325
August 2018

Part No: E69325

Copyright © 2000, 2018, Oracle and/or its affiliates. All rights reserved.

This software and related documentation are provided under a license agreement containing restrictions on use and disclosure and are protected by intellectual property laws. Except as expressly permitted in your license agreement or allowed by law, you may not use, copy, reproduce, translate, broadcast, modify, license, transmit, distribute, exhibit, perform, publish, or display any part, in any form, or by any means. Reverse engineering, disassembly, or decompilation of this software, unless required by law for interoperability, is prohibited.

The information contained herein is subject to change without notice and is not warranted to be error-free. If you find any errors, please report them to us in writing.

If this is software or related documentation that is delivered to the U.S. Government or anyone licensing it on behalf of the U.S. Government, then the following notice is applicable:

U.S. GOVERNMENT END USERS: Oracle programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, delivered to U.S. Government end users are "commercial computer software" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, use, duplication, disclosure, modification, and adaptation of the programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, shall be subject to license terms and license restrictions applicable to the programs. No other rights are granted to the U.S. Government.

This software or hardware is developed for general use in a variety of information management applications. It is not developed or intended for use in any inherently dangerous applications, including applications that may create a risk of personal injury. If you use this software or hardware in dangerous applications, then you shall be responsible to take all appropriate fail-safe, backup, redundancy, and other measures to ensure its safe use. Oracle Corporation and its affiliates disclaim any liability for any damages caused by use of this software or hardware in dangerous applications.

Oracle and Java are registered trademarks of Oracle and/or its affiliates. Other names may be trademarks of their respective owners.

Intel and Intel Xeon are trademarks or registered trademarks of Intel Corporation. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. AMD, Opteron, the AMD logo, and the AMD Opteron logo are trademarks or registered trademarks of Advanced Micro Devices. UNIX is a registered trademark of The Open Group.

This software or hardware and documentation may provide access to or information about content, products, and services from third parties. Oracle Corporation and its affiliates are not responsible for and expressly disclaim all warranties of any kind with respect to third-party content, products, and services unless otherwise set forth in an applicable agreement between you and Oracle. Oracle Corporation and its affiliates will not be responsible for any loss, costs, or damages incurred due to your access to or use of third-party content, products, or services, except as set forth in an applicable agreement between you and Oracle.

Access to Oracle Support

Oracle customers that have purchased support have access to electronic support through My Oracle Support. For information, visit <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=info> or visit <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=trs> if you are hearing impaired.

Référence: E69325

Copyright © 2000, 2018, Oracle et/ou ses affiliés. Tous droits réservés.

Ce logiciel et la documentation qui l'accompagne sont protégés par les lois sur la propriété intellectuelle. Ils sont concédés sous licence et soumis à des restrictions d'utilisation et de divulgation. Sauf stipulation expresse de votre contrat de licence ou de la loi, vous ne pouvez pas copier, reproduire, traduire, diffuser, modifier, accorder de licence, transmettre, distribuer, exposer, exécuter, publier ou afficher le logiciel, même partiellement, sous quelque forme et par quelque procédé que ce soit. Par ailleurs, il est interdit de procéder à toute ingénierie inverse du logiciel, de le désassembler ou de le décompiler, excepté à des fins d'interopérabilité avec des logiciels tiers ou tel que prescrit par la loi.

Les informations fournies dans ce document sont susceptibles de modification sans préavis. Par ailleurs, Oracle Corporation ne garantit pas qu'elles soient exemptes d'erreurs et vous invite, le cas échéant, à lui en faire part par écrit.

Si ce logiciel, ou la documentation qui l'accompagne, est livré sous licence au Gouvernement des Etats-Unis, ou à quiconque qui aurait souscrit la licence de ce logiciel pour le compte du Gouvernement des Etats-Unis, la notice suivante s'applique :

U.S. GOVERNMENT END USERS: Oracle programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, delivered to U.S. Government end users are "commercial computer software" pursuant to the applicable Federal Acquisition Regulation and agency-specific supplemental regulations. As such, use, duplication, disclosure, modification, and adaptation of the programs, including any operating system, integrated software, any programs installed on the hardware, and/or documentation, shall be subject to license terms and license restrictions applicable to the programs. No other rights are granted to the U.S. Government.

Ce logiciel ou matériel a été développé pour un usage général dans le cadre d'applications de gestion des informations. Ce logiciel ou matériel n'est pas conçu ni n'est destiné à être utilisé dans des applications à risque, notamment dans des applications pouvant causer un risque de dommages corporels. Si vous utilisez ce logiciel ou ce matériel dans le cadre d'applications dangereuses, il est de votre responsabilité de prendre toutes les mesures de secours, de sauvegarde, de redondance et autres mesures nécessaires à son utilisation dans des conditions optimales de sécurité. Oracle Corporation et ses affiliés déclinent toute responsabilité quant aux dommages causés par l'utilisation de ce logiciel ou matériel pour des applications dangereuses.

Oracle et Java sont des marques déposées d'Oracle Corporation et/ou de ses affiliés. Tout autre nom mentionné peut correspondre à des marques appartenant à d'autres propriétaires qu'Oracle.

Intel et Intel Xeon sont des marques ou des marques déposées d'Intel Corporation. Toutes les marques SPARC sont utilisées sous licence et sont des marques ou des marques déposées de SPARC International, Inc. AMD, Opteron, le logo AMD et le logo AMD Opteron sont des marques ou des marques déposées d'Advanced Micro Devices. UNIX est une marque déposée de The Open Group.

Ce logiciel ou matériel et la documentation qui l'accompagne peuvent fournir des informations ou des liens donnant accès à des contenus, des produits et des services émanant de tiers. Oracle Corporation et ses affiliés déclinent toute responsabilité ou garantie expresse quant aux contenus, produits ou services émanant de tiers, sauf mention contraire stipulée dans un contrat entre vous et Oracle. En aucun cas, Oracle Corporation et ses affiliés ne sauraient être tenus pour responsables des pertes subies, des coûts occasionnés ou des dommages causés par l'accès à des contenus, produits ou services tiers, ou à leur utilisation, sauf mention contraire stipulée dans un contrat entre vous et Oracle.

Accès aux services de support Oracle

Les clients Oracle qui ont souscrit un contrat de support ont accès au support électronique via My Oracle Support. Pour plus d'informations, visitez le site <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=info> ou le site <http://www.oracle.com/pls/topic/lookup?ctx=acc&id=trs> si vous êtes malentendant.

Contents

Using This Documentation	9
1 Introduction to Oracle Solaris Cluster Hardware	11
Installing Oracle Solaris Cluster Hardware	11
Supported Configurations	12
▼ Installing Oracle Solaris Cluster Hardware	13
Maintaining Oracle Solaris Cluster Hardware	14
Powering Oracle Solaris Cluster Hardware On and Off	15
Dynamic Reconfiguration Operations for Oracle Solaris Cluster Nodes	15
▼ Dynamic Reconfiguration Operations in a Cluster With Dynamic Reconfiguration-Enabled Servers	15
Local and Multihost Disks in an Oracle Solaris Cluster Environment	16
Removable Media in an Oracle Solaris Cluster Environment	16
SAN Solutions in an Oracle Solaris Cluster Environment	17
Hardware Restrictions	17
2 Installing Cluster Interconnect Hardware and Configuring VLANs	19
Interconnect Requirements and Restrictions	19
Cluster Interconnect and Routing	20
Cluster Interconnect Speed Requirements	20
Ethernet Switch Configuration When in the Cluster Interconnect	20
Requirements When Using Jumbo Frames	20
Requirements and Restrictions When Using Sun InfiniBand from Oracle in the Cluster Interconnect	21
Requirements for Socket Direct Protocol Over an Oracle Solaris Cluster Interconnect	22
Installing Ethernet or InfiniBand Cluster Interconnect Hardware	22
▼ How to Install Ethernet or InfiniBand Transport Cables and Transport Junctions	23

Configuring VLANs as Private Interconnect Networks	24
Configuring SR-IOV Devices with VLANs as Private Interconnect Networks	27
3 Maintaining Cluster Interconnect Hardware	29
Maintaining Interconnect Hardware in a Running Cluster	29
▼ How to Add an Interconnect Component	30
▼ How to Replace an Interconnect Component	31
▼ How to Remove an Interconnect Component	33
▼ How to Upgrade Transport Adapter Firmware	35
4 Installing and Maintaining Public Network Hardware	37
Public Network Hardware: Requirements When Using Jumbo Frames	37
Installing Public Network Hardware	38
Installing Public Network Hardware: Where to Go From Here	38
Maintaining Public Network Hardware in a Running Cluster	38
Adding Public Network Adapters	39
Replacing Public Network Adapters	39
Removing Public Network Adapters	40
5 Maintaining Platform Hardware	41
Mirroring Internal Disks on Servers that Use Internal Hardware Disk Mirroring or Integrated Mirroring	41
▼ How to Configure Internal Disk Mirroring After the Cluster Is Established	42
▼ How to Remove an Internal Disk Mirror	44
Configuring Cluster Nodes With a Single, Dual-Port HBA	45
Risks and Trade-offs When Using One Dual-Port HBA	46
Supported Configurations When Using a Single, Dual-Port HBA	47
Cluster Configuration When Using Solaris Volume Manager and a Single Dual- Port HBA	47
Cluster Configuration When Using Solaris Volume Manager for Sun Cluster and a Single Dual-Port HBA	48
Kernel Cage Dynamic Reconfiguration Recovery	49
Preparing the Cluster for Kernel Cage Dynamic Reconfiguration	49
▼ How to Recover From an Interrupted Kernel Cage Dynamic Reconfiguration Operation	50

6 Campus Clustering With Oracle Solaris Cluster Software	53
Requirements for Designing a Campus Cluster	53
Selecting Networking Technologies	54
Connecting to Storage	54
Shared Data Storage	55
Complying With Quorum Device Requirements	55
Replicating Solaris Volume Manager Disksets	55
Guidelines for Designing a Campus Cluster	56
Determining the Number of Rooms in Your Cluster	56
Deciding How to Use Quorum Devices	60
Determining Campus Cluster Connection Technologies	63
Cluster Interconnect Technologies	63
Storage Area Network Technologies	64
Installing and Configuring Interconnect, Storage, and Fibre Channel Hardware	65
Calculating Buffer Credits	65
Additional Campus Cluster Configuration Examples	66
7 Verifying Oracle Solaris Cluster Hardware Redundancy	71
Testing Node Redundancy	72
▼ How to Test Device Group Redundancy Using Resource Group Failover	72
Testing Cluster Interconnect Redundancy	73
▼ How to Test Cluster Interconnects	73
Testing Public Network Redundancy	74
▼ How to Test Public Network Redundancy	74
Index	77

Using This Documentation

- **Overview** – Describes how to install and administer basic Oracle Solaris Cluster hardware components
- **Audience** – Technicians, system administrators, and authorized service providers
- **Required knowledge** – Advanced experience troubleshooting and replacing hardware

Product Documentation Library

Documentation and resources for this product and related products are available at http://docs.oracle.com/cd/E69294_01.

Feedback

Provide feedback about this documentation at <http://www.oracle.com/goto/docfeedback>.

◆◆◆ CHAPTER 1

Introduction to Oracle Solaris Cluster Hardware

The information and procedures in this book apply to Oracle Solaris Cluster 4.0 and subsequent releases.

This chapter provides overview information on cluster hardware. The chapter also provides overviews of the tasks that are involved in installing and maintaining this hardware specifically in an Oracle Solaris Cluster environment.

This chapter contains the following information:

- [“Installing Oracle Solaris Cluster Hardware” on page 11](#)
- [“Maintaining Oracle Solaris Cluster Hardware” on page 14](#)
- [“Powering Oracle Solaris Cluster Hardware On and Off” on page 15](#) [“Powering Oracle Solaris Cluster Hardware On and Off” on page 15](#)
- [“Dynamic Reconfiguration Operations for Oracle Solaris Cluster Nodes” on page 15](#)
- [“Local and Multihost Disks in an Oracle Solaris Cluster Environment” on page 16](#)
- [“Removable Media in an Oracle Solaris Cluster Environment” on page 16](#)
- [“SAN Solutions in an Oracle Solaris Cluster Environment” on page 17](#)
- [“Hardware Restrictions” on page 17](#)

Installing Oracle Solaris Cluster Hardware

The following procedure lists the tasks for installing a cluster and where to find instructions.

TABLE 1 Task Map: Installing Cluster Hardware

Task	For Instructions
Plan for cluster hardware capacity, space, and power requirements.	The site planning documentation that shipped with your nodes and other hardware
Verify that your system is supported.	“Supported Configurations” on page 12
Install the nodes.	The documentation that shipped with your nodes

Task	For Instructions
Install the administrative console.	The documentation that shipped with your administrative console
Install a console access device.	The documentation that shipped with your hardware
Install the cluster interconnect hardware.	Chapter 2, “Installing Cluster Interconnect Hardware and Configuring VLANs”
Install the public network hardware.	Chapter 4, “Installing and Maintaining Public Network Hardware”
Install and configure the shared disk storage arrays.	The Oracle Solaris Cluster guide that pertains to your storage device as well as to the device's own documentation
Install the Oracle Solaris Operating System and Oracle Solaris Cluster software.	Chapter 2, “Installing Software on Global-Cluster Nodes” in <i>Installing and Configuring an Oracle Solaris Cluster 4.4 Environment</i>
Configure the cluster interconnects.	Chapter 3, “Establishing the Global Cluster” in <i>Installing and Configuring an Oracle Solaris Cluster 4.4 Environment</i>

Supported Configurations

Depending on your platform, Oracle Solaris Cluster software supports the following configurations:

- **SPARC:** Oracle Solaris Cluster software supports from one to 16 cluster nodes in a cluster. Different hardware configurations impose additional limits on the maximum number of nodes that you can configure in a cluster composed of SPARC-based systems. See [“Oracle Solaris Cluster Topologies”](#) in *Concepts for Oracle Solaris Cluster 4.4* for the supported configurations.
- **x86:** Oracle Solaris Cluster software supports from one to eight cluster nodes in a cluster. Different hardware configurations impose additional limits on the maximum number of nodes that you can configure in a cluster composed of x86-based systems. See [“Oracle Solaris Cluster Topologies”](#) in *Concepts for Oracle Solaris Cluster 4.4* for the supported configurations.

All nodes in the cluster must have the same architecture. All nodes in the cluster must run the same version of the Oracle Solaris OS. Nodes in the same cluster must have the same OS and architecture, as well as similar processing, memory, and I/O capability, to enable failover to occur without significant degradation in performance. Because of the possibility of failover, every node must have enough excess capacity to support the workload of all nodes for which they are a backup or secondary.

Cluster nodes are generally attached to one or more multihost storage devices. Nodes that are not attached to multihost devices can use a cluster file system to access the data on multihost

devices. For example, one scalable services configuration enables nodes to service requests without being directly attached to multihost devices.

In addition, nodes in parallel database configurations share concurrent access to all the disks.

- See [“Multihost Devices”](#) in *Concepts for Oracle Solaris Cluster 4.4* for information about concurrent access to disks.
- See [“Clustered Pair Topology”](#) in *Concepts for Oracle Solaris Cluster 4.4* and [“Clustered Pair Topology”](#) in *Concepts for Oracle Solaris Cluster 4.4* for more information about parallel database configurations and scalable topology.

Public network adapters attach nodes to the public networks, providing client access to the cluster.

Cluster members communicate with the other nodes in the cluster through one or more physically independent networks. This set of physically independent networks is referred to as the *cluster interconnect*.

Every node in the cluster is aware when another node joins or leaves the cluster. Additionally, every node in the cluster is aware of the resources that are running locally as well as the resources that are running on the other cluster nodes.

▼ Installing Oracle Solaris Cluster Hardware

1. Plan for cluster hardware capacity, space, and power requirements.

For more information, see the site planning documentation that shipped with your servers and other hardware. See [“Hardware Restrictions”](#) on [page 17](#) for critical information about hardware restrictions with Oracle Solaris Cluster.

2. Install the nodes.

For server installation instructions, see the documentation that shipped with your servers.

3. Install the administrative console.

For more information, see the documentation that shipped with your administrative console.

4. Install a console access device.

Use the procedure that is indicated for your type of console access device.

5. Install the cluster interconnect and public network hardware.

For installation instructions, see [Chapter 2, “Installing Cluster Interconnect Hardware and Configuring VLANs”](#).

6. Install and configure the storage arrays.

Perform the service procedures that are indicated for your type of storage hardware.

7. Install the Oracle Solaris Operating System and Oracle Solaris Cluster software.

For more information, see the [Installing and Configuring an Oracle Solaris Cluster 4.4 Environment](#).

8. Plan, install, and configure resource groups and data services.

For more information, see [Planning and Administering Data Services for Oracle Solaris Cluster 4.4](#).

Maintaining Oracle Solaris Cluster Hardware

This guide augments documentation that ships with your hardware components by providing information on maintaining the hardware *specifically in an Oracle Solaris Cluster environment*. [Table 2, “Sample Differences Between Servicing Standalone and Cluster Hardware,”](#) on [page 14](#) describes some of the differences between maintaining cluster hardware and maintaining standalone hardware.

TABLE 2 Sample Differences Between Servicing Standalone and Cluster Hardware

Task	Standalone Hardware	Cluster Hardware
Shutting down a node	Use the shutdown command.	To perform an orderly node shutdown, first use the <code>clnode evacuate</code> to switch device groups and resource groups to another node. Then shut down the node by running the <code>shutdown(8)</code> command.
Adding a disk	Perform a reconfiguration boot or use <code>devfsadm</code> to assign a logical device name to the disk. You also need to run volume manager commands to configure the new disk if the disks are under volume management control.	Use the <code>devfsadm</code> , <code>cldevice populate</code> , and <code>cldevice</code> or <code>scdidadm</code> commands. You also need to run volume manager commands to configure the new disk if the disks are under volume management control.
Adding a transport adapter or public network adapter (PNA)	Perform an orderly node shutdown, then install the public network adapter. After you install the network adapter, update the <code>/etc/hostname.adapter</code> and <code>/etc/inet/hosts</code> files.	Perform an orderly node shutdown, then install the public network adapter. After you install the public network adapter, update the <code>/etc/hostname.adapter</code> and <code>/etc/inet/hosts</code> files. Finally, add this public network adapter to an IPMP group.

Powering Oracle Solaris Cluster Hardware On and Off

Consider the following when powering on and powering off cluster hardware.

- Use shut down and boot procedures in the [Administering an Oracle Solaris Cluster 4.4 Configuration](#) for nodes in a running cluster.
- Use the power-on and power-off procedures in the manuals that shipped with the hardware *only* for systems that are newly installed or are in the process of being installed.



Caution - After the cluster is online and a user application is accessing data on the cluster, do not use the power-on and power-off procedures listed in the manuals that came with the hardware.

Dynamic Reconfiguration Operations for Oracle Solaris Cluster Nodes

The Oracle Solaris Cluster environment supports Oracle Solaris dynamic reconfiguration operations on qualified servers. Contact your service provider for a list of storage arrays that are qualified for use with servers that are enabled with dynamic reconfiguration.

Note - Review the documentation for the Oracle Solaris dynamic reconfiguration feature on your hardware platform *before* you use the dynamic reconfiguration feature with Oracle Solaris Cluster software. All of the requirements, procedures, and restrictions that are documented for the Oracle Solaris dynamic reconfiguration feature also apply to Oracle Solaris Cluster dynamic reconfiguration support (except for the operating environment quiescence operation).

▼ Dynamic Reconfiguration Operations in a Cluster With Dynamic Reconfiguration-Enabled Servers

Some Oracle Solaris Cluster procedures instruct the user to shut down and power off a cluster node before you add, remove, or replace a transport adapter or a public network adapter (PNA).

However, if the node is a server that is enabled with the dynamic reconfiguration feature, the user does *not* have to power off the node before you add, remove, or replace the transport adapter or PNA. Instead, do the following:

For conceptual information about Oracle Solaris Cluster support of the dynamic reconfiguration feature, see [“Dynamic Reconfiguration Support” in *Concepts for Oracle Solaris Cluster 4.4*](#).

1. **Follow the steps in [“How to Remove Cluster Transport Cables, Transport Adapters, and Transport Switches” in *Administering an Oracle Solaris Cluster 4.4 Configuration*](#), including steps for disabling and removing the transport adapter or PNA from the active cluster interconnect.**
2. **Skip any step that instructs you to power off the node, where the purpose of the power-off is to add, remove, or replace a transport adapter or PNA.**
3. **Perform the dynamic reconfiguration operation (add, remove, or replace) on the transport adapter or PNA.**

Local and Multihost Disks in an Oracle Solaris Cluster Environment

Two sets of storage arrays reside within a cluster: local disks and multihost disks.

- Local disks are directly connected to a single node and hold the Oracle Solaris Operating System and other nonshared data.
- Multihost disks are connected to more than one node and hold client application data and other files that need to be accessed from multiple nodes.

For more conceptual information on multihost disks and local disks, see the [Concepts for Oracle Solaris Cluster 4.4](#).

Removable Media in an Oracle Solaris Cluster Environment

Removable media include tape and CD-ROM drives, which are local devices. This guide does not contain procedures for adding, removing, or replacing removable media as highly available storage arrays. Although tape and CD-ROM drives are global devices, these drives are not supported as highly available. Thus, this guide focuses on disk drives as global devices.

Although tape and CD-ROM drives are not supported as highly available in a cluster environment, you can access tape and CD-ROM drives that are not local to your system. All the various density extensions (such as h, b, l, n, and u) are mapped so that the tape drive can be accessed from any node in the cluster.

Install, remove, replace, and use tape and CD-ROM drives as you would in a noncluster environment. For procedures about how to install, remove, and replace tape and CD-ROM drives, see the documentation that shipped with your hardware.

SAN Solutions in an Oracle Solaris Cluster Environment

You cannot have a single point of failure in a SAN configuration that is in an Oracle Solaris Cluster environment. For information about how to install and configure a SAN configuration, see your SAN documentation.

Hardware Restrictions

The following restrictions apply to hardware in all Oracle Solaris Cluster configurations.

- All nodes in the same cluster must be of the same architecture. They must be all SPARC-based systems or all x86-based systems.
- Multihost tape, CD-ROM, and DVD-ROM are not supported.
- Alternate pathing (AP) is not supported.
- Storage devices with more than a single path from a given cluster node to the enclosure are not supported except for the following storage devices:
 - Oracle's Sun StorEdge™ A3500, for which two paths are supported to each of two nodes.
 - Devices using Oracle Solaris I/O multipathing, formerly Sun StorEdge Traffic Manager.
 - EMC storage devices that use EMC PowerPath software.
 - Oracle's Sun StorEdge 9900 storage devices that use HDLM.
- SunVTS™ software is not supported.

Installing Cluster Interconnect Hardware and Configuring VLANs

This chapter describes the procedures to install cluster interconnect hardware. Where appropriate, this chapter includes separate procedures for the interconnects that Oracle Solaris Cluster software supports:

- Ethernet
- InfiniBand

This chapter contains the following information:

- [“Installing Ethernet or InfiniBand Cluster Interconnect Hardware” on page 22](#)
- [“Configuring VLANs as Private Interconnect Networks” on page 24](#)

Use the following information to learn more about cluster interconnects:

- For conceptual information about cluster interconnects, see [“Cluster Interconnect” in *Concepts for Oracle Solaris Cluster 4.4*](#).
- For information about how to administer cluster interconnects, see [Chapter 7, “Administering Cluster Interconnects and Public Networks” in *Administering an Oracle Solaris Cluster 4.4 Configuration*](#).
- For information about how to configure VLANs as private interconnects, see [“Cluster Interconnect” in *Installing and Configuring an Oracle Solaris Cluster 4.4 Environment*](#).

Interconnect Requirements and Restrictions

This section contains requirements on interconnect operation when using certain special features.

Cluster Interconnect and Routing

Heartbeat packets that are sent over the cluster interconnect are not IP based. As a result, these packets cannot be routed. If you install a router between two cluster nodes that are connected through cluster interconnects, heartbeat packets cannot find their destination. Your cluster consequently fails to work correctly.

To ensure that your cluster works correctly, you must set up the cluster interconnect in the same layer 2 (data link) network and in the same broadcast domain. The cluster interconnect must be located in the same layer 2 network and broadcast domain even if the cluster nodes are located in different, remote data centers. Cluster nodes that are arranged remotely are described in more detail in [Chapter 6, “Campus Clustering With Oracle Solaris Cluster Software”](#).

Cluster Interconnect Speed Requirements

An interconnect path is one network step in the cluster private network: from a node to a node, from a node to a switch, or from the switch to another node. Each path in your cluster interconnect must use the same networking technology.

All interconnect paths must also operate at the same speed. This means, for example, that if you are using Ethernet components that are capable of operating at different speeds, and if your cluster configuration does not allow these components to automatically negotiate a common network speed, you must configure them to operate at the same speed.

Ethernet Switch Configuration When in the Cluster Interconnect

When configuring Ethernet switches for your cluster private interconnect, disable the spanning tree algorithm on ports that are used for the interconnect.

Requirements When Using Jumbo Frames

If you use Scalable Data Services and jumbo frames on your public network, ensure that the Maximum Transfer Unit (MTU) of the private network is the same size or larger than the MTU of your public network.

Note - Scalable services cannot forward public network packets that are larger than the MTU size of the private network. The scalable services application instances will not receive those packets.

Consider the following information when configuring jumbo frames:

- The maximum MTU size for an InfiniBand interface is typically less than the maximum MTU size for an Ethernet interface.
- If you use switches in your private network, ensure they are configured to the MTU sizes of the private network interfaces.



Caution - If the switches are not configured to the MTU sizes of the private network interfaces, the cluster interconnect might not stay online.

For information about how to configure jumbo frames, see the documentation that shipped with your network interface card. See your Oracle Solaris OS documentation or contact your Oracle sales representative for other Oracle Solaris restrictions.

Requirements and Restrictions When Using Sun InfiniBand from Oracle in the Cluster Interconnect

The following requirements and guidelines apply to Oracle Solaris Cluster configurations that use Sun InfiniBand adapters from Oracle:

- A two-node cluster must use InfiniBand switches. You cannot directly connect the InfiniBand adapters to each other.
- If only one InfiniBand adapter is installed on a cluster node, each of its two ports must be connected to a different InfiniBand switch.
- If two InfiniBand adapters are installed in a cluster node, leave the second port on each adapter unused for interconnect purposes. For example, connect port 1 on HCA 1 to switch 1 and connect port 1 on HCA 2 to switch 2 when using these connections as a cluster interconnect.

Requirements for Socket Direct Protocol Over an Oracle Solaris Cluster Interconnect

In an Oracle Solaris Cluster configuration that uses an InfiniBand interconnect, applications can use Socket Direct Protocol (SDP) by configuring SDP to use the `clprivnetN` network device. If there is a failure at the port of the HCA or switch, Automatic Path Migration (APM) fails over all live SDP sessions to the standby HCA port in a manner that is transparent to the application. APM is a built-in failover facility that is included in the InfiniBand software.

APM cannot be performed if the standby port is connected to a different switch partition, and the application must explicitly reestablish SDP sessions to recover. To ensure that APM can be performed successfully, observe the following requirements:

- If redundant InfiniBand switches are set up as a cluster interconnect, you must use multiple HCAs. Both ports of an HCA must be connected to the same switch, and only one of the two HCA ports can be configured as a cluster interconnect.
- If only one InfiniBand switch is set up as a cluster interconnect, you can use only one HCA. Both ports of the HCA must be connected to the same InfiniBand partition on the switch, and both ports can be configured as a cluster interconnect.

Installing Ethernet or InfiniBand Cluster Interconnect Hardware

The following table lists procedures for installing Ethernet or InfiniBand cluster interconnect hardware. Perform the procedures in the order that they are listed. This section contains the procedure for installing cluster hardware during an *initial installation* of a cluster, before you install Oracle Solaris Cluster software.

TABLE 3 Installing Ethernet Cluster Interconnect Hardware

Task	For Instructions
Install the transport adapters.	The documentation that shipped with your nodes and host adapters
Install the transport cables.	“How to Install Ethernet or InfiniBand Transport Cables and Transport Junctions” on page 23
If your cluster contains more than two nodes, install a transport junction (switch).	“How to Install Ethernet or InfiniBand Transport Cables and Transport Junctions” on page 23

▼ How to Install Ethernet or InfiniBand Transport Cables and Transport Junctions

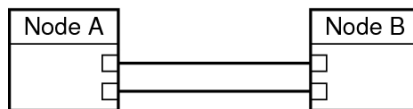
Use this procedure to install Ethernet or InfiniBand transport cables and transport junctions (switches).

1. **If not already installed, install transport adapters in your cluster nodes.**
See the documentation that shipped with your host adapters and node hardware.
2. **If necessary, install transport junctions and optionally configure the transport junctions' IP addresses.**

Note - (InfiniBand Only) If you install one InfiniBand adapter on a cluster node, two InfiniBand switches are required. Each of the two ports must be connected to a different InfiniBand switch.

If two InfiniBand adapters are connected to a cluster node, use only one port on each adapter for the interconnect and have it connected to an InfiniBand switch. The second port of the adapter can be connected but must not be used as an interconnect. Do not connect ports of the two InfiniBand adapters to the same InfiniBand switch.

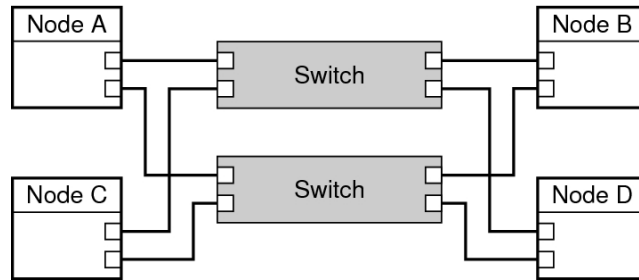
3. **Install the transport cables.**
 - **(Ethernet Only)** As the following figure shows, a cluster with only two nodes can use a point-to-point connection, requiring no transport junctions.



(Ethernet Only) For a point-to-point connection, you can use either UTP or fibre. With fibre, use a standard patch cable. A crossover cable is unnecessary. With UTP, see your network interface card documentation to determine whether you need a crossover cable.

Note - (Ethernet Only) You can optionally use transport junctions in a two-node cluster. If you use a transport junction in a two-node cluster, you can more easily add additional nodes later. To ensure redundancy and availability, always use two transport junctions.

- As the following figure shows, a cluster with more than two nodes requires transport junctions. These transport junctions are Ethernet or InfiniBand switches (customer-supplied).



See Also To install and configure the Oracle Solaris Cluster software with the new interconnect, see [Chapter 2, “Installing Software on Global-Cluster Nodes”](#) in *Installing and Configuring an Oracle Solaris Cluster 4.4 Environment*.

(Ethernet Only) To configure jumbo frames on the interconnect, review the requirements in [“Requirements When Using Jumbo Frames”](#) on page 20.

Configuring VLANs as Private Interconnect Networks

Oracle Solaris Cluster software supports the use of private interconnect networks over switch-based virtual local area networks (VLANs). In a switch-based VLAN environment, Oracle Solaris Cluster software enables multiple clusters and nonclustered systems to share an Ethernet transport junction (switch) in two different configurations.

Note - Even if clusters share the same switch, create a separate VLAN for each cluster.

By default, Oracle Solaris Cluster uses the same set of IP addresses on the private interconnect. Creating a separate VLAN for each cluster ensures that IP traffic from one cluster does not interfere with IP traffic from another cluster. Unless you have customized the default IP address for the private interconnect, as described in [“How to Change the Private Network Address or Address Range of an Existing Cluster”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*, create a separate VLAN for each cluster.

The implementation of switch-based VLAN environments is vendor-specific. Because each switch manufacturer implements VLAN differently, the following guidelines address Oracle Solaris Cluster software requirements with regard to configuring VLANs with cluster interconnects.

- You must understand your capacity needs before you set up a VLAN configuration. You must know the minimum bandwidth necessary for your interconnect and application traffic. For the best results, set the Quality of Service (QOS) level for each VLAN to accommodate basic cluster traffic and the desired application traffic. Ensure that the bandwidth that is allocated to each VLAN extends from node to node.

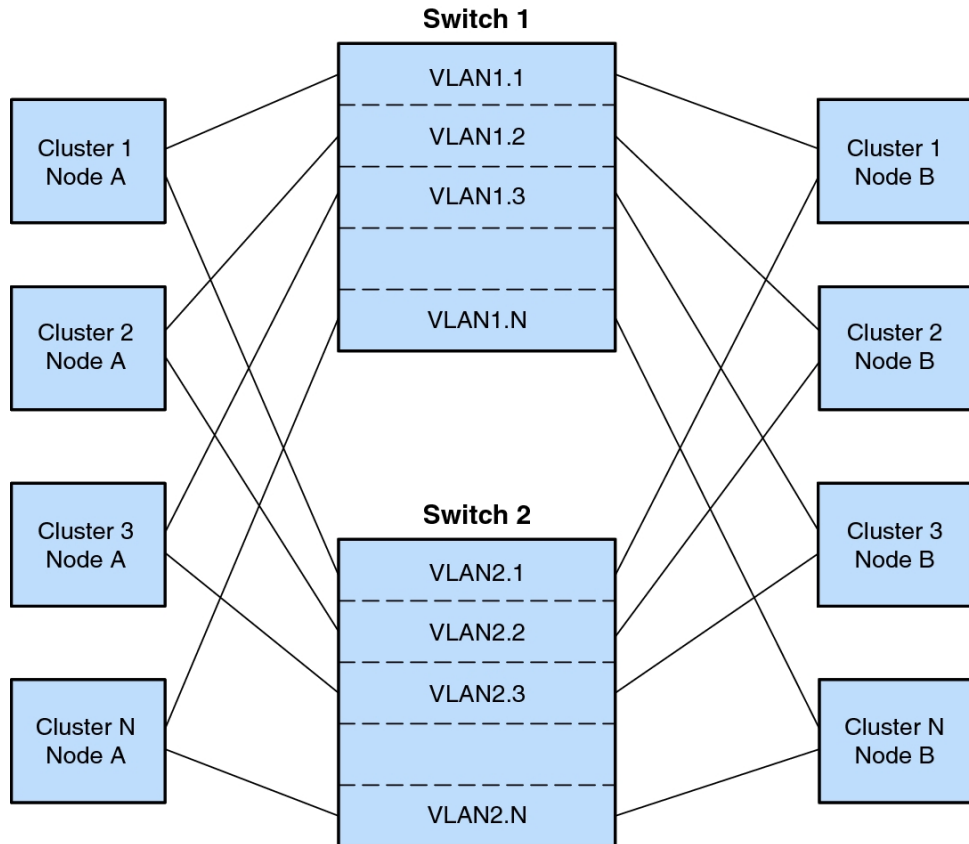
To determine the basic cluster traffic requirements, use the following equation. In this equation, n equals the number of nodes in the configuration, and s equals the number of switches per VLAN.

$$n (s-1) \times 10\text{Mb}$$

- The use of two cluster interconnects provides higher availability than one interconnect. If the number of available adapter ports is limited, you can use tagged VLANs to share the same adapter with both the private and public network. For more information, see the guidelines for tagged VLAN adapters in [“Transport Adapters” in *Installing and Configuring an Oracle Solaris Cluster 4.4 Environment*](#).
- Interconnect traffic must be placed in the highest-priority queue.
- All ports must be equally serviced, similar to a round robin or first-in, first-out model.
- You must verify that you have correctly configured your VLANs to prevent path timeouts.

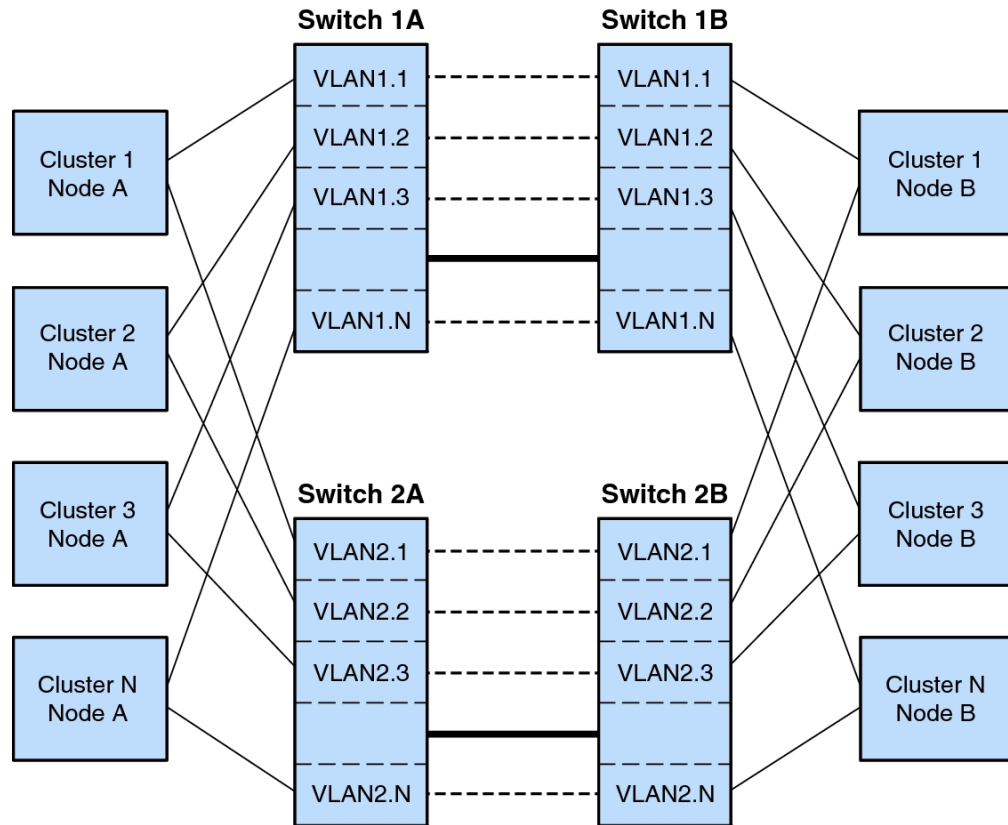
The first VLAN configuration enables nodes from multiple clusters to send interconnect traffic across one pair of Ethernet transport junctions. Oracle Solaris Cluster software requires a minimum of one transport junction, and each transport junction must be part of a VLAN that is located on a different switch. The following figure is an example of the first VLAN configuration in a two-node cluster. VLAN configurations are not limited to two-node clusters.

FIGURE 1 First VLAN Configuration



The second VLAN configuration uses the same transport junctions for the interconnect traffic of multiple clusters. However, the second VLAN configuration has two pairs of transport junctions that are connected by links. This configuration enables VLANs to be supported in a campus cluster configuration with the same restrictions as other campus cluster configurations. The following figure illustrates the second VLAN configuration.

FIGURE 2 Second VLAN Configuration



Configuring SR-IOV Devices with VLANs as Private Interconnect Networks

Oracle Solaris Cluster supports SR-IOV (VF) devices for public network and transport.

For transport usage, the traffic will be seen by all VF devices and hence different guest domain clusters have to choose a unique default network to avoid conflict. To use a default network without any conflict, set the `pvId` property with a tagged-vlan ID created on the network switch.

```
# ldm set-io pvid=<tagged-vlanid> VF-Device
```

For example:

```
# ldm set-io pvid=2859 /SYS/MB/NET2/IOVNET.PF1.VF0
# ldm ls-io -l
```

```
/SYS/MB/NET2/IOVNET.PF1.VF0          VF      pci_2    guest1
[pci@300/pci@3/network@0,81]
  Class properties [NETWORK]
    mac-addr = 00:14:4f:fb:73:51
    port-vlan-id = 2859
    mtu = 1500
```

Maintaining Cluster Interconnect Hardware

This chapter describes the procedures to maintain cluster interconnect hardware. The procedures in this chapter apply to all interconnects that Oracle Solaris Cluster software supports:

- Ethernet
- InfiniBand

This chapter contains the following procedures:

- [“How to Add an Interconnect Component” on page 30](#)
- [“How to Replace an Interconnect Component” on page 31](#)
- [“How to Remove an Interconnect Component” on page 33](#)
- [“How to Upgrade Transport Adapter Firmware” on page 35](#)

For more information, see the following documentation:

- For conceptual information about cluster interconnects, see [“Cluster Interconnect” in *Concepts for Oracle Solaris Cluster 4.4*](#).
- For information about administering cluster interconnects, see [“Administering the Cluster Interconnects” in *Administering an Oracle Solaris Cluster 4.4 Configuration*](#).

Maintaining Interconnect Hardware in a Running Cluster

The following table lists procedures about maintaining cluster interconnect hardware.

TABLE 4 Task Map: Maintaining Cluster Interconnect Hardware

Task	Instructions
Add an interconnect component.	“How to Add an Interconnect Component” on page 30
Replace an interconnect component.	“How to Replace an Interconnect Component” on page 31
Remove an interconnect component.	“How to Remove an Interconnect Component” on page 33
Upgrade transport adapter firmware	“How to Upgrade Transport Adapter Firmware” on page 35

Interconnect components include the following components:

- Transport adapter
- Transport cable
- Transport junction (switch)

▼ How to Add an Interconnect Component

This procedure defines interconnect component as any one of the following components:

- Transport adapter
- Transport cable
- Transport junction (switch)

This section contains the procedure for adding interconnect components to nodes in a running cluster.

Before You Begin This procedure relies on the following prerequisites and assumptions:

- Your cluster is operational and all nodes are powered on.
- If virtual local area networks (VLANs) are configured, more than one cluster might be impacted by removing a transport junction. Ensure that all clusters are prepared for the removal of a transport junction. Also, record the configuration information of the transport junction you plan to replace and configure the new transport junction accordingly.

For more information about how to configure VLANs, see [“Configuring VLANs as Private Interconnect Networks” on page 24](#).

You can also add an interconnect cable, switch, or private adapter using the Oracle Solaris Cluster Manager GUI. For GUI log-in instructions, see [“How to Access Oracle Solaris Cluster Manager” in *Administering an Oracle Solaris Cluster 4.4 Configuration*](#).

- 1. Determine if you need to shut down and power off the node that is to be connected to the interconnect component you are adding.**
 - If you are adding a transport junction, you do not need to shut down and power off the node. Proceed to [Step 2](#).
 - If you are adding a transport cable, you do not need to shut down and power off the node. Proceed to [Step 2](#).
 - If your node has dynamic reconfiguration enabled and you are replacing a transport adapter, you do not need to shut down and power off the node. Proceed to [Step 2](#).

- If your node does *not* have dynamic reconfiguration enabled and you are adding a transport adapter, shut down and power off the node with the transport adapter you are adding.

For the full procedure about shutting down a node, see [Chapter 3, “Shutting Down and Booting a Cluster”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.

2. Install the interconnect component.

- If you are using an Ethernet or InfiniBand interconnect, see [“How to Install Ethernet or InfiniBand Transport Cables and Transport Junctions”](#) on page 23 for cabling diagrams and considerations.
- For the procedure about installing transport adapters or setting transport adapter DIP switches, see the documentation that shipped with your host adapter and node hardware.
- If your interconnect uses jumbo frames, review the requirements in [“Requirements When Using Jumbo Frames”](#) on page 20.

3. If you shut down the node in [Step 1](#), perform a reconfiguration boot to update the new Oracle Solaris device files and links. Otherwise, skip this step.

- See Also
- To reconfigure Oracle Solaris Cluster software with the new interconnect component, see [Chapter 7, “Administering Cluster Interconnects and Public Networks”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.

▼ How to Replace an Interconnect Component

This procedure defines interconnect component as any one of the following components:

- Transport adapter
- Transport cable
- Transport junction (switch)



Caution - You must maintain at least one cluster interconnect between the nodes of a cluster. The cluster does not function without a working cluster interconnect. You can check the status of the interconnect with the `clinterconnect status` command.

For more details about checking the status of the cluster interconnect, see [“How to Check the Status of the Cluster Interconnect”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.

You might perform this procedure in the following scenarios:

- You need to replace a failed transport adapter.

- You need to replace a failed transport cable.
- You need to replace a failed transport junction.

For conceptual information about transport adapters, transport cables, and transport junction, see [“Cluster Interconnect” in *Concepts for Oracle Solaris Cluster 4.4*](#).

Before You Begin This procedure relies on the following prerequisites and assumptions.

- Your cluster has another functional interconnect path to maintain cluster communications while you perform this procedure.
- Your cluster is operational and all nodes are powered on.
- Identify the interconnect component that you want to replace. Remove that interconnect component from the cluster configuration by using the procedure in [“How to Remove Cluster Transport Cables, Transport Adapters, and Transport Switches” in *Administering an Oracle Solaris Cluster 4.4 Configuration*](#).
- If virtual local area networks (VLANs) are configured, more than one cluster might be impacted by removing a transport junction. Ensure that all clusters are prepared for the removal of a transport junction. Also, record the configuration information of the transport junction you plan to replace and configure the new transport junction accordingly.

For more information about how to configure VLANs, see [“Configuring VLANs as Private Interconnect Networks” on page 24](#).

1. Determine if you need to shut down and power off the node that is connected to the interconnect component you are replacing.

- If you are replacing a transport junction, you do not need to shut down and power off the node. Proceed to [Step 2](#).
- If you are replacing a transport cable, you do not need to shut down and power off the node. Proceed to [Step 2](#).
- If your node has dynamic reconfiguration enabled and you are replacing a transport adapter, you do not need to shut down and power off the node. Proceed to [Step 2](#).
- If your node does *not* have dynamic reconfiguration enabled and you are replacing a transport adapter, shut down and power off the node with the transport adapter you are replacing.

For the full procedure about how to shut down a node, see [Chapter 3, “Shutting Down and Booting a Cluster” in *Administering an Oracle Solaris Cluster 4.4 Configuration*](#).

2. Disconnect the failed interconnect component from other cluster devices.

For the procedure about how to disconnect cables from transport adapters, see the documentation that shipped with your host adapter and node.

3. Connect the new interconnect component to other cluster devices.

- If you are replacing an Ethernet or InfiniBand interconnect, see [“How to Install Ethernet or InfiniBand Transport Cables and Transport Junctions”](#) on page 23 for cabling diagrams and considerations.
 - If your interconnect uses jumbo frames, review the requirements in [“Requirements When Using Jumbo Frames”](#) on page 20.
4. **If you shut down the node in [Step 1](#), perform a reconfiguration boot to update the new Oracle Solaris device files and links. Otherwise, skip this step.**

See Also To reconfigure Oracle Solaris Cluster software with the new interconnect component, see [“How to Add Cluster Transport Cables, Transport Adapters, or Transport Switches”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.

▼ How to Remove an Interconnect Component

This procedure defines interconnect component as any one of the following components:

- Transport adapter
- Transport cable
- Transport junction (switch)



Caution - You must maintain at least one cluster interconnect between the nodes of a cluster. The cluster does not function without a working cluster interconnect. You can check the status of the interconnect with the `cLinterconnect status` command.

For more details about checking the status of the cluster interconnect, see [“How to Check the Status of the Cluster Interconnect”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.

You might perform this procedure in the following scenarios:

- You need to remove an unused transport adapter.
- You need to remove an unused transport cable.
- You need to remove an unused transport junction.
- You want to migrate from a two-node cluster that uses switches to a point-to-point configuration.

For conceptual information about transport adapters, transport cables, and transport junctions, see [“Cluster Interconnect”](#) in *Concepts for Oracle Solaris Cluster 4.4*.

Before You Begin This procedure assumes that your cluster is operational and all nodes are powered on. Before you perform this procedure, perform the following tasks:

- If you are migrating from a two-node cluster that uses switches to a point-to-point configuration, install a crossover cable before you remove a switch.
- Identify the interconnect component that you want to remove. Remove that interconnect component from the cluster configuration by using the procedure in [“How to Remove Cluster Transport Cables, Transport Adapters, and Transport Switches”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.
- If you plan to use virtual local area networks (VLANs) in your cluster interconnect, configure the transport junction. For more information about how to configure VLANs, see [“Configuring VLANs as Private Interconnect Networks”](#) on page 24.

1. Determine if you need to shut down and power off the node that is connected to the interconnect component you are removing.

- If you are removing a transport junction, you do not need to shut down and power off the node. Proceed to [Step 2](#).
- If you are removing a transport cable, you do not need to shut down and power off the node. Proceed to [Step 2](#).
- If your node has dynamic reconfiguration enabled and you are removing a transport adapter, you do not need to shut down and power off the node. Proceed to [Step 2](#).
- If your node does *not* have dynamic reconfiguration enabled and you are removing a transport adapter, shut down and power off the node with the transport adapter you are removing.

For the full procedure about shutting down a node, see [Chapter 3, “Shutting Down and Booting a Cluster”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.

2. Disconnect the interconnect component from other cluster devices.

For instructions on how to disconnect cables from transport adapters, see the documentation that shipped with your host adapter and node.

3. Remove the interconnect component.

For instructions on how to remove the interconnect component, see the documentation that shipped with your host adapter, nodes, or switch.

4. If you shut down the node in [Step 1](#), perform a reconfiguration boot to update the new Oracle Solaris device files and links. Otherwise, skip this step.

See Also To reconfigure Oracle Solaris Cluster software with the new interconnect component, see [“How to Add Cluster Transport Cables, Transport Adapters, or Transport Switches”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.

▼ How to Upgrade Transport Adapter Firmware

You might perform this procedure in the following scenarios:

- You want to use firmware bug fixes.
- You want to use new firmware features.

Use this procedure to update transport adapter firmware.

Before You Begin To perform this procedure, become an administrator that provides `solaris.cluster.read` and `solaris.cluster.modify` authorization.

1. **Determine the resource groups and the device groups that are online on the node. This node is the node where you are upgrading transport adapter firmware.**

Use the following command:

```
# clresourcegroup status -n nodename
# cldevicegroup status -n nodename
```

Note the device groups, the resource groups, and the node list for the resource groups. You will need this information to restore the cluster to its original configuration in [Step 4](#).

2. **Migrate the resource groups and device groups off the node on which you plan to upgrade the firmware.**

```
# clnode evacuate fromnode
```

3. **Perform the firmware upgrade.**

This process might require you to boot into noncluster mode. If it does, boot the node into cluster mode before proceeding. For the procedure about how to upgrade your transport adapter firmware, see the patch documentation.

4. **If you moved device groups off their original node in [Step 2](#), restore the device groups that you identified in [Step 1](#) to their original node.**

Perform the following step for each device group you want to return to the original node.

```
# cldevicegroup switch -n nodename devicegroup1[ devicegroup2 ...]
```

`-n nodename` The node to which you are restoring device groups.

`devicegroup1[
devicegroup2 ...]` The device group or groups that you are restoring to the node.

In these commands, *devicegroup* is one or more device groups that are returned to the node.

5. If you moved resource groups off their original node in Step 2 restore the resource groups that you identified in Step 1 to their original node.

Perform the following step for each resource group you want to return to the original node.

```
# clresourcegroup switch -n nodename resourcegroup1[ resourcegroup2 ...]
```

nodename For failover resource groups, the node to which the groups are returned.
For scalable resource groups, the node list to which the groups are returned.

resourcegroup1[resourcegroup2 ...] The resource group or groups that you are returning to the node or nodes.

resourcegroup The resource group that is returned to the node or nodes.

Installing and Maintaining Public Network Hardware

This chapter contains information about how to maintain public network hardware. This chapter covers the following topics.

- “Public Network Hardware: Requirements When Using Jumbo Frames” on page 37
- “Installing Public Network Hardware” on page 38
- “Maintaining Public Network Hardware in a Running Cluster” on page 38

For conceptual information on cluster interconnects and public network interfaces, see your [Concepts for Oracle Solaris Cluster 4.4](#).

For information on how to administer public network interfaces, see your [Administering an Oracle Solaris Cluster 4.4 Configuration](#).

Note - Some hardware drivers are no longer available in the Oracle Solaris 11 OS. These include drivers for Sun Gigabit Ethernet and Sun GigaSwift Ethernet. For up-to-date lists, see the [Oracle Solaris Hardware Compatibility Lists](#).

Public Network Hardware: Requirements When Using Jumbo Frames

If you use Scalable Data Services and jumbo frames on your public network, ensure that the Maximum Transfer Unit (MTU) of the private network is the same size or larger than the MTU of your public network.

Note - Scalable services cannot forward public network packets that are larger than the MTU size of the private network. The scalable services application instances will not receive those packets.

Consider the following information when configuring jumbo frames:

- The maximum MTU size for an InfiniBand interface is typically less than the maximum MTU size for an Ethernet interface.
- If you use switches in your private network, ensure they are configured to the MTU sizes of the private network interfaces.

For information about how to configure jumbo frames, see the documentation that shipped with your network interface card. See your Oracle Solaris OS documentation or contact your Oracle sales representative for other Oracle Solaris restrictions.

Installing Public Network Hardware

This section covers installing cluster hardware during an *initial cluster installation*, before *Oracle Solaris Cluster software is installed*.

Physically installing public network adapters to a node in a cluster is no different from adding public network adapters in a noncluster environment.

For the procedure about how to add public network adapters, see the documentation that shipped with your nodes and public network adapters.

Installing Public Network Hardware: Where to Go From Here

Install the cluster software and configure the public network hardware after you have installed all other hardware. To review the task map about how to install cluster hardware, see [“Installing Oracle Solaris Cluster Hardware” on page 11](#).

If your network uses jumbo frames, review the requirements in [“Public Network Hardware: Requirements When Using Jumbo Frames” on page 37](#).

Maintaining Public Network Hardware in a Running Cluster

The following table lists procedures about how to maintain public network hardware.

TABLE 5 Task Map: Maintaining Public Network Hardware

Task	Information
Add public network adapters.	“Adding Public Network Adapters” on page 39
Replace public network adapters.	“Replacing Public Network Adapters” on page 39
Remove public network adapters.	“Removing Public Network Adapters” on page 40

Adding Public Network Adapters

Physically adding public network adapters to a node in a cluster is no different from adding public network adapters in a noncluster environment. For the procedure about how to add public network adapters, see the hardware documentation that shipped with your node and public network adapters.

Once the adapters are physically installed, Oracle Solaris Cluster requires that they be configured in an IPMP group.

If your network uses jumbo frames, review the requirements in [“Public Network Hardware: Requirements When Using Jumbo Frames” on page 37](#) and see the documentation that shipped with your network interface card for information about how to configure jumbo frames.

Adding Public Network Adapters: Where to Go From Here

To add a new public network adapter to an IPMP group, see [Chapter 1, “Administering TCP/IP Networks” in *Administering TCP/IP Networks, IPMP, and IP Tunnels in Oracle Solaris 11.4*](#).

Replacing Public Network Adapters

For cluster-specific commands and guidelines about how to replace public network adapters, see [“Administering the Public Network” in *Administering an Oracle Solaris Cluster 4.4 Configuration*](#).

For procedures about how to administer public network connections, see the [Chapter 2, “About IPMP Administration” in *Administering TCP/IP Networks, IPMP, and IP Tunnels in Oracle Solaris 11.4*](#).

For the procedure about removing public network adapters, see the hardware documentation that shipped with your node and public network adapters.

Replacing Public Network Adapters: Where to Go From Here

To add the new public network adapter to a IPMP group, see the [Administering an Oracle Solaris Cluster 4.4 Configuration](#).

Removing Public Network Adapters

For cluster-specific commands and guidelines about how to remove public network adapters, see the [Administering an Oracle Solaris Cluster 4.4 Configuration](#).

For procedures about how to administer public network connections, see the [Chapter 2, “About IPMP Administration” in Administering TCP/IP Networks, IPMP, and IP Tunnels in Oracle Solaris 11.4](#).

For the procedure about how to remove public network adapters, see the hardware documentation that shipped with your node and public network adapters.

Maintaining Platform Hardware

This chapter contains information about node hardware in a cluster environment. It contains the following topics:

- [“Mirroring Internal Disks on Servers that Use Internal Hardware Disk Mirroring or Integrated Mirroring”](#) on page 41
- [“Configuring Cluster Nodes With a Single, Dual-Port HBA”](#) on page 45
- [“Kernel Cage Dynamic Reconfiguration Recovery”](#) on page 49

Mirroring Internal Disks on Servers that Use Internal Hardware Disk Mirroring or Integrated Mirroring

Some servers support the mirroring of internal hard drives (internal hardware disk mirroring or integrated mirroring) to provide redundancy for node data. To use this feature in a cluster environment, follow the steps in this section.

The best way to set up hardware disk mirroring is to perform RAID configuration during cluster installation, before you configure multipathing. For instructions on performing this configuration, see the *[Installing and Configuring an Oracle Solaris Cluster 4.4 Environment](#)*. If you need to change your mirroring configuration after you have established the cluster, you must perform some cluster-specific steps to clean up the device IDs, as described in the procedure that follows.

Note - Specific servers might have additional restrictions. See the documentation that shipped with your server hardware.

For specifics about how to configure your server's internal disk mirroring, refer to the documents that shipped with your server and the [raidctl\(8\)](#) man page.

▼ How to Configure Internal Disk Mirroring After the Cluster Is Established

Before You Begin This procedure assumes that you have already installed your hardware and software and have established the cluster. To configure an internal disk mirror during cluster installation, see the [Installing and Configuring an Oracle Solaris Cluster 4.4 Environment](#).



Caution - If there are state database replicas on the disk that you are mirroring, you must recreate them during this procedure.

1. If necessary, prepare the node for establishing the mirror.

a. Determine the resource groups and device groups that are running on the node.

Record this information because you use it later in this procedure to return resource groups and device groups to the node.

Use the following command:

```
# clresourcegroup status -n nodename
# cldevicegroup status -n nodename
```

b. If necessary, move all resource groups and device groups off the node.

```
# clnode evacuate fromnode
```

2. Configure the internal mirror.

```
# raidctl -c clt0d0 clt1d0
```

`-c clt0d0 clt1d0` Creates the mirror of primary disk to the mirror disk. Enter the name of your primary disk as the first argument. Enter the name of the mirror disk as the second argument.

3. Boot the node into single user mode.

```
# shutdown -g 0 -i 0 -y
ok> boot -s
```

4. Clean up the device IDs.

Use the following command:

```
# cldevice repair /dev/rdisk/clt0d0
```

`/dev/rdisk/clt0d0` Updates the cluster's record of the device IDs for the primary disk. Enter the name of your primary disk as the argument.

- 5. Confirm that the mirror has been created and only the primary disk is visible to the cluster.**

```
# cldevice list
```

The command lists only the primary disk, and not the mirror disk, as visible to the cluster.

- 6. Boot the node back into cluster mode.**

```
# shutdown -g 0 -i 0 -y
ok> boot -s
```

- 7. If you are using Solaris Volume Manager and if the state database replicas are on the primary disk, recreate the state database replicas.**

```
# metadb -a /dev/rdisk/clt0d0s4
```

- 8. If you moved device groups off the node in [Step 1](#), restore device groups to the original node.**

Perform the following step for each device group you want to return to the original node.

```
# cldevicegroup switch -n nodename devicegroup1[ devicegroup2 ...]
```

`-n nodename` The node to which you are restoring device groups.

`devicegroup1[devicegroup2 ...]` The device group or groups that you are restoring to the node.

- 9. If you moved resource groups off the node in [Step 1](#), move all resource groups back to the node.**

Perform the following step for each resource group you want to return to the original node.

```
# clresourcegroup switch -n nodename resourcegroup1[ resourcegroup2 ...]
```

`nodename` For failover resource groups, the node to which the groups are returned. For scalable resource groups, the node list to which the groups are returned.

`resourcegroup1[resourcegroup2 ...]` The resource group or groups that you are returning to the node or nodes.

▼ How to Remove an Internal Disk Mirror

1. If necessary, prepare the node for removing the mirror.

a. Determine the resource groups and device groups that are running on the node.

Record this information because you use this information later in this procedure to return resource groups and device groups to the node.

Use the following command:

```
# clresourcegroup status -n nodename
# cldevicegroup status -n nodename
```

b. If necessary, move all resource groups and device groups off the node.

```
# clnode evacuate fromnode
```

2. Remove the internal mirror.

```
# raidctl -d clt0d0
```

`-d clt0d0` Deletes the mirror of primary disk to the mirror disk. Enter the name of your primary disk as the argument.

3. Boot the node into single user mode.

```
# reboot -- -S
```

4. Clean up the device IDs.

Use the following command:

```
# cldevice repair /dev/rdisk/clt0d0 /dev/rdisk/clt1d0
```

`/dev/rdisk/
clt0d0 /dev/rdisk/
clt1d0` Updates the cluster's record of the device IDs. Enter the names of your disks separated by spaces.

5. Confirm that the mirror has been deleted and that both disks are visible.

```
# cldevice list
```

The command lists both disks as visible to the cluster.

6. Boot the node back into cluster mode.

```
# shutdown -g 0 -i 0 -y
ok> boot
```

7. If you are using Solaris Volume Manager and if the state database replicas are on the primary disk, recreate the state database replicas.

```
# metadb -c 3 -ag /dev/rdisk/clt0d0s4
```

8. If you moved device groups off the node in [Step 1](#), restore the device groups to the original node.

```
# cldevicegroup switch -n nodename devicegroup1 devicegroup2 ...
```

-n nodename The node where you are restoring device groups.

devicegroup1[The device group or groups that you are restoring to the node.
devicegroup2 ...]

9. If you moved resource groups off the node in [Step 1](#), restore the resource groups and device groups to the original node.

Perform the following step for each resource group you want to return to the original node.

```
# clresourcegroup switch -n nodename resourcegroup[ resourcegroup2 ...]
```

nodename For failover resource groups, the node to which the groups are restored.
For scalable resource groups, the node list to which the groups are restored.

resourcegroup[The resource group or groups that you are restoring to the node or nodes.
resourcegroup2 ...]

Configuring Cluster Nodes With a Single, Dual-Port HBA

This section explains the use of dual-port host bus adapters (HBAs) to provide both connections to shared storage in the cluster. While Oracle Solaris Cluster supports this configuration, it is less redundant than the recommended configuration. You *must* understand the risks that a dual-port HBA configuration poses to the availability of your application, if you choose to use this configuration.

This section contains the following topics:

- [“Risks and Trade-offs When Using One Dual-Port HBA” on page 46](#)

- [“Supported Configurations When Using a Single, Dual-Port HBA” on page 47](#)
- [“Cluster Configuration When Using Solaris Volume Manager and a Single Dual-Port HBA” on page 47](#)
- [“Cluster Configuration When Using Solaris Volume Manager for Sun Cluster and a Single Dual-Port HBA” on page 48](#)

Risks and Trade-offs When Using One Dual-Port HBA

You should strive for as much separation and hardware redundancy as possible when connecting each cluster node to shared data storage. This approach provides the following advantages to your cluster:

- The best assurance of high availability for your clustered application
- Good failure isolation
- Good maintenance robustness

Oracle Solaris Cluster is usually layered on top of a volume manager, mirrored data with independent I/O paths, or a multipathed I/O link to a hardware RAID arrangement. Therefore, the cluster software does not expect a node ever to ever lose access to shared data. These redundant paths to storage ensure that the cluster can survive any single failure.

Oracle Solaris Cluster does support certain configurations that use a single, dual-port HBA to provide the required two paths to the shared data. However, using a single, dual-port HBA for connecting to shared data increases the vulnerability of your cluster. If this single HBA fails and takes down both ports connected to the storage device, the node is unable to reach the stored data. How the cluster handles such a dual-port failure depends on several factors:

- The cluster configuration
- The volume manager configuration
- The node on which the failure occurs
- The state of the cluster when the failure occurs

If you choose one of these configurations for your cluster, you must understand that the supported configurations mitigate the risks to high availability and the other advantages. The supported configurations do not eliminate these previously mentioned risks.

Supported Configurations When Using a Single, Dual-Port HBA

Oracle Solaris Cluster supports the following volume manager configurations when you use a single, dual-port HBA for connecting to shared data:

- Solaris Volume Manager with more than one disk in each diskset and no dual-string mediators configured. For details about this configuration, see [“Cluster Configuration When Using Solaris Volume Manager and a Single Dual-Port HBA” on page 47](#).
- Solaris Volume Manager for Sun Cluster. For details about this configuration, see [“Cluster Configuration When Using Solaris Volume Manager for Sun Cluster and a Single Dual-Port HBA” on page 48](#).

Cluster Configuration When Using Solaris Volume Manager and a Single Dual-Port HBA

If the Solaris Volume Manager metadbs lose replica quorum for a diskset on a cluster node, the volume manager panics the cluster node. Oracle Solaris Cluster then takes over the diskset on a surviving node and your application fails over to a secondary node.

To ensure that the node panics and is fenced off if it loses its connection to shared storage, configure each metaset with at least two disks. In this configuration, the metadbs stored on the disks create their own replica quorum for each diskset.

Dual-string mediators are not supported in Solaris Volume Manager configurations that use a single dual-port HBA. Using dual-string mediators prevents the service from failing over to a new node.

Configuration Requirements

When configuring Solaris Volume Manager metasets, ensure that each metaset contains at least two disks. Do not configure dual-string mediators.

Expected Failure Behavior with Solaris Volume Manager

When a dual-port HBA fails with both ports in this configuration, the cluster behavior depends on whether the affected node is primary for the diskset.

- If the affected node is primary for the diskset, Solaris Volume Manager panics that node because it lacks required state database replicas. Your cluster reforms with the nodes that achieve quorum and brings the diskset online on a new primary node.
- If the affected node is not primary for the diskset, your cluster remains in a degraded state.

Failure Recovery with Solaris Volume Manager

Follow the instructions for replacing an HBA in your storage device documentation.

Cluster Configuration When Using Solaris Volume Manager for Sun Cluster and a Single Dual-Port HBA

Because Solaris Volume Manager for Sun Cluster uses raw disks only and is specific to Oracle Real Application Clusters (RAC), no special configuration is required.

Expected Failure Behavior with Solaris Volume Manager for Sun Cluster

When a dual-port HBA fails and takes down both ports in this configuration, the cluster behavior depends on whether the affected node is the current master for the multi-owner diskset.

- If the affected node is the current master for the multi-owner diskset, the node does not panic. If any other node fails or is rebooted, the affected node will panic when it tries to update the replicas. The volume manager chooses a new master for the diskset if the surviving nodes can achieve quorum.
- If the affected node is not the current master for the multi-owner diskset, the node remains up but the device group is in a degraded state. If an additional failure affects the master node and Solaris Volume Manager for Sun Cluster attempts to remaster the diskset on the node with the failed paths, that node will also panic. A new master will be chosen if any surviving nodes can achieve quorum.

Failure Recovery with Solaris Volume Manager for Sun Cluster

Follow the instructions for replacing an HBA in your storage device documentation.

Kernel Cage Dynamic Reconfiguration Recovery

When you perform a dynamic reconfiguration remove operation on a memory board with kernel cage memory, the affected node becomes unresponsive so heartbeat monitoring for that node is suspended on all other nodes and the node's quorum vote count is set to 0. After dynamic reconfiguration is completed, the heartbeat monitoring of the affected node is automatically re-enabled and the quorum vote count is reset to 1. If the dynamic reconfiguration operation does not complete, you might need to manually recover. For general information about dynamic reconfiguration, see [“Dynamic Reconfiguration Support” in *Concepts for Oracle Solaris Cluster 4.4*](#).

The `monitor-heartbeat` subcommand is not supported in an exclusive-IP zone cluster. For more information about this command, see the `cluster(8CL)` man page.

Preparing the Cluster for Kernel Cage Dynamic Reconfiguration

When you use a dynamic reconfiguration operation to remove a system board containing kernel cage memory (memory used by the Oracle Solaris OS), the system must be quiesced in order to allow the memory contents to be copied to another system board. In a clustered system, the tight coupling between cluster nodes means that the quiescing of one node for repair can cause operations on non-quiesced nodes to be delayed until the repair operation is complete and the node is unquiesced. For this reason, using dynamic reconfiguration to remove a system board containing kernel cage memory from a cluster node requires careful planning and preparation.

Use the following information to reduce the impact of the dynamic reconfiguration quiesce on the rest of the cluster:

- I/O operations for file systems or global device groups with their primary or secondary on the quiesced node will hang until the node is unquiesced. If possible, ensure that the node being repaired is not the primary for any global file systems or device groups.
- I/O to SVM multi-owner disksets that include the quiesced node will hang until the node is unquiesced.

- Updates to the CCR require communication between all cluster members. Any operations that result in CCR updates should not be performed while the dynamic reconfiguration operation is ongoing. Configuration changes are the most common cause of CCR updates.
- Many cluster commands result in communication among cluster nodes. Refrain from running cluster commands during the dynamic reconfiguration operation.
- Applications and cluster resources on the node being quiesced will be unavailable for the duration of the dynamic reconfiguration event. The time required to move applications and resources to another node should be weighed against the expected outage time of the dynamic reconfiguration event.
- Scalable applications such as Oracle RAC often have a different membership standard, and have communication and synchronization actions among members. Scalable application instances on the node to be repaired should be brought offline before you initiate the dynamic reconfiguration operation.

▼ How to Recover From an Interrupted Kernel Cage Dynamic Reconfiguration Operation

If the dynamic reconfiguration operation does not complete, perform the following steps to re-enable heartbeat timeout monitoring for that node and to reset the quorum vote count.

- 1. If dynamic reconfiguration does not complete successfully, manually re-enable heartbeat timeout monitoring.**

From a single cluster node (which is not the node where the dynamic reconfiguration operation was performed), run the following command.

```
# cluster monitor-heartbeat
```

Use this command only in the global zone. Messages display indicating that monitoring has been enabled.

- 2. If the node that was dynamically reconfigured paused during boot, allow it to finish booting and join the cluster membership.**

If the node is at the ok prompt, boot it now.

- 3. Verify that the node is now part of the cluster membership and check the quorum vote count of the cluster nodes by running the following command on a single node in the cluster.**

```
# clquorum status
--- Quorum Votes by Node (current status) ---
```

Node Name	Present	Possible	Status
-----	-----	-----	-----
pnode1	1	1	Online
pnode2	1	1	Online
pnode3	0	0	Online

- If one of the nodes has a vote count of 0, reset its vote count to 1 by running the following command on a single node in the cluster.**

```
# clquorum votecount -n nodename 1
```

nodename The hostname of the node that has a quorum vote count of 0.

- Verify that all nodes now have a quorum vote count of 1.**

```
# clquorum status
--- Quorum Votes by Node (current status) ---
```

Node Name	Present	Possible	Status
-----	-----	-----	-----
pnode1	1	1	Online
pnode2	1	1	Online
pnode3	1	1	Online

Campus Clustering With Oracle Solaris Cluster Software

In campus clustering, nodes or groups of nodes are located in separate rooms, sometimes several kilometers apart. In addition to providing the usual benefits of using an Oracle Solaris Cluster, correctly designed campus clusters can generally survive the loss of any single room and continue to provide their services.

This chapter introduces the basic concepts of campus clustering and provides some configuration and setup examples. The following topics are covered:

- “Requirements for Designing a Campus Cluster” on page 53
- “Guidelines for Designing a Campus Cluster” on page 56
- “Determining Campus Cluster Connection Technologies” on page 63
- “Installing and Configuring Interconnect, Storage, and Fibre Channel Hardware” on page 65
- “Additional Campus Cluster Configuration Examples” on page 66

This chapter does not explain clustering, provide information about clustering administration, or furnish details about hardware installation and configuration. For conceptual and administrative information, see the *Concepts for Oracle Solaris Cluster 4.4* and the *Administering an Oracle Solaris Cluster 4.4 Configuration*, respectively.

Requirements for Designing a Campus Cluster

When designing your campus cluster, all of the requirements for a standard cluster still apply. Plan your cluster to eliminate any single point of failure in nodes, cluster interconnect, data storage, and public network. Just as in the standard cluster, a campus cluster requires redundant connections and switches. Disk multipathing helps ensure that each node can access each shared storage device. These concerns are universal for Oracle Solaris Cluster.

After you have a valid cluster plan, follow the requirements in this section to ensure a correct campus cluster. To achieve maximum benefits from your campus cluster, consider implementing the [“Guidelines for Designing a Campus Cluster” on page 56](#).

Note - This chapter describes ways to design your campus cluster using fully tested and supported hardware components and transport technologies. You can also design your campus cluster according to Oracle Solaris Cluster's specification, regardless of the components used.

To build a specifications-based campus cluster, contact your Oracle representative, who will assist you with the design and implementation of your specific configuration. This process ensures that the configuration that you implement complies with the specification guidelines, is interoperable, and is supportable.

Selecting Networking Technologies

Your campus cluster must observe all requirements and limitations of the technologies that you choose to use. [“Determining Campus Cluster Connection Technologies” on page 63](#) provides a list of tested technologies and their known limitations.

When planning your cluster interconnect, remember that campus clustering requires redundant network connections.

Connecting to Storage

A campus cluster must include at least two rooms using two independent SANs to connect to the shared storage. See [Figure 3, “Basic Three-Room, Two-Node Campus Cluster Configuration With Multipathing,” on page 58](#) for an illustration of this configuration.

If you are using Oracle Real Application Clusters (RAC), all nodes that support Oracle RAC must be fully connected to the shared storage devices. Also, all rooms of a specifications-based campus cluster must be fully connected to the shared storage devices.

See [“Quorum in Clusters With Four Rooms or More” on page 61](#) for a description of a campus cluster with both direct and indirect storage connections.

Shared Data Storage

Your campus cluster must use SAN-supported storage devices for shared storage. When planning the cluster, ensure that it adheres to the SAN requirements for all storage connections. See the [SAN Storage documentation site \(http://www.oracle.com/technetwork/server-storage/san-storage/documentation/index.html\)](http://www.oracle.com/technetwork/server-storage/san-storage/documentation/index.html) for information about SAN requirements.

Oracle Solaris Cluster software supports storage-based replication, which moves the work of data replication off the cluster nodes and onto the storage device. Storage-based data replication can simplify the infrastructure required, which can be useful in campus cluster configurations..

For more information on data replication and supported software, see [Chapter 4, “Data Replication Approaches”](#) in *Administering an Oracle Solaris Cluster 4.4 Configuration*.

Complying With Quorum Device Requirements

You must use a quorum device for a two-node cluster. For larger clusters, a quorum device is optional. These are standard cluster requirements.

In addition, you can configure quorum devices to ensure that specific rooms can form a cluster in the event of a failure. For guidelines about where to locate your quorum device, see [“Deciding How to Use Quorum Devices”](#) on page 60.

Replicating Solaris Volume Manager Disksets

If you use Solaris Volume Manager as your volume manager for shared device groups, carefully plan the distribution of your replicas. In two-room configurations, all disksets should be configured with an additional replica in the room that houses the cluster quorum device.

For example, in three-room two-node configurations, a single room houses both the quorum device and at least one extra disk that is configured in each of the disksets. Each diskset should have extra replicas in the third room.

Note - You can use a quorum disk for these replicas.

Refer to your Solaris Volume Manager documentation for details about configuring diskset replicas.

Guidelines for Designing a Campus Cluster

In planning a campus cluster, your goal is to build a cluster that can at least survive the loss of a room and continue to provide services. The concept of a room must shape your planning of redundant connectivity, storage replication, and quorum. Use the following guidelines to assist in managing these design considerations.

Determining the Number of Rooms in Your Cluster

The concept of a room, or location, adds a layer of complexity to the task of designing a campus cluster. Think of a *room* as a functionally independent hardware grouping, such as a node and its attendant storage, or a quorum device that is physically separated from any nodes. Each room is separated from other rooms to increase the likelihood of failover and redundancy in case of accident or failure. The definition of a room therefore depends on the type of failure to safeguard against, as described in the following table.

TABLE 6 Definitions of "Room"

Failure Scenario	Sample Definitions of "Room"
Power-line failure	Isolated and independent power supplies
Minor accidents, furniture collapse, water seepage	Different parts of a physical room
Small fire, fire sprinklers starting	Different physical areas (for example, sprinkler zone)
Structural failure, building-wide fire	Different buildings
Large-scale natural disaster (for example, earthquake or flood)	Different corporate campuses up to several kilometers apart

Oracle Solaris Cluster does support two-room campus clusters. These clusters are valid and might offer nominal insurance against disasters. However, consider adding a small third room, possibly even a secure closet or vault (with a separate power supply and correct cabling), to contain the quorum device or a third server.

Whenever a two-room campus cluster loses a room, it has only a 50 percent chance of remaining available. If the room with fewest quorum votes is the surviving room, the surviving nodes cannot form a cluster. In this case, your cluster requires manual intervention from your Oracle service provider before it can become available.

The advantage of a three-room or larger cluster is that, if any one of the three rooms is lost, automatic failover can be achieved. Only a correctly configured three-room or larger campus cluster can guarantee system availability if an entire room is lost (assuming no other failures).

Three-Room Campus Cluster Examples

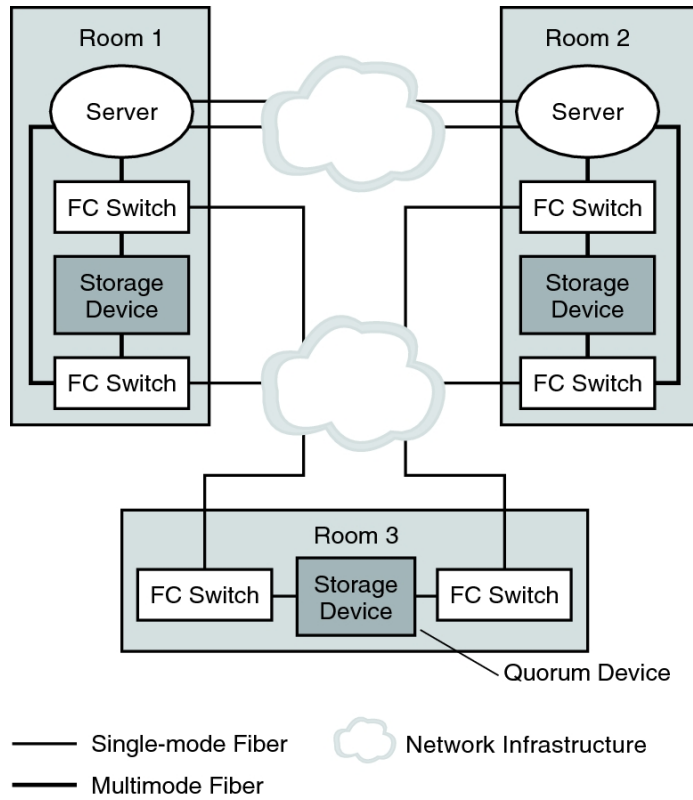
A three-room campus cluster configuration supports up to eight nodes. Three rooms enable you to arrange your nodes and quorum device so that your campus cluster can reliably survive the loss of a single room and still provide cluster services. The following example configurations all follow the campus cluster requirements and the design guidelines described in this chapter.

- [Figure 3, “Basic Three-Room, Two-Node Campus Cluster Configuration With Multipathing,” on page 58](#) shows a three-room, two-node campus cluster. In this arrangement, two rooms each contain a single node and an equal number of disk arrays to mirror shared data. The third room contains at least one disk subsystem, attached to both nodes and configured with a quorum device.
- [Figure 4, “Minimum Three-Room, Two-Node Campus Cluster Configuration Without Multipathing,” on page 59](#) shows an alternative three-room, two-node campus cluster.
- [Figure 5, “Three-Room, Three-Node Campus Cluster Configuration,” on page 60](#) shows a three-room, three-node cluster. In this arrangement, two rooms each contain one node and an equal number of disk arrays. The third room contains a small server, which eliminates the need for a storage array to be configured as a quorum device.

Note - These examples illustrate general configurations and are not intended to indicate required or recommended setups. For simplicity, the diagrams and explanations concentrate only on features that are unique to understanding campus clustering. For example, public-network Ethernet connections are not shown.

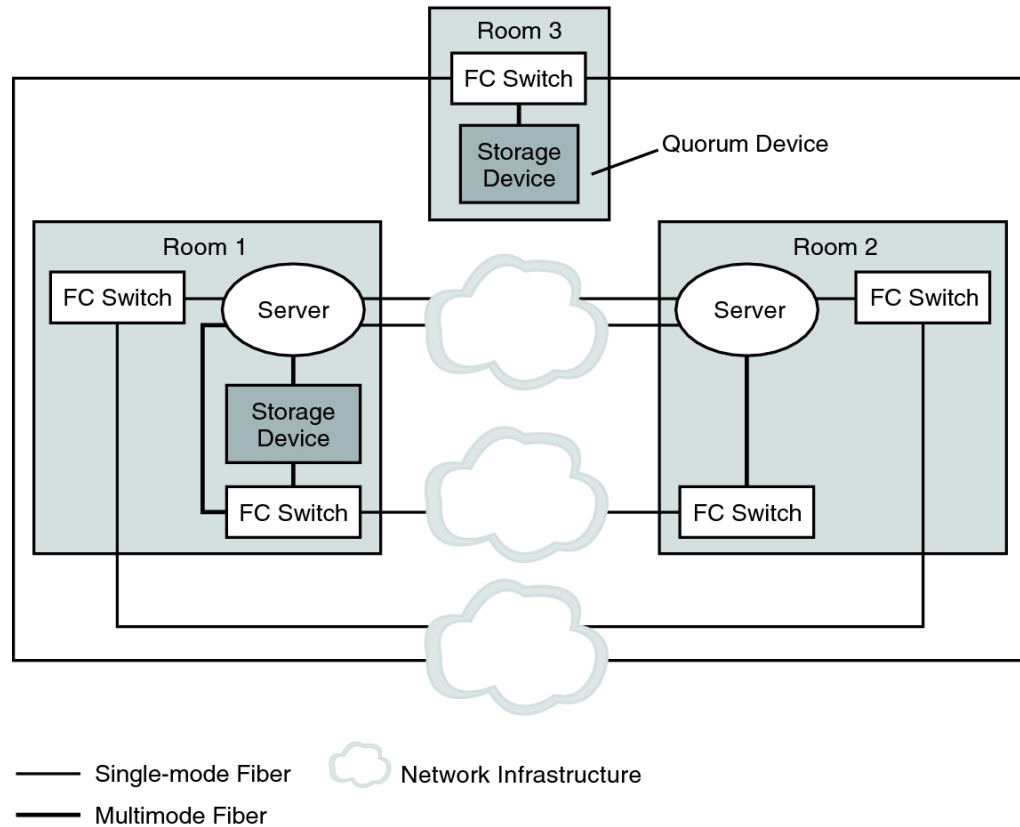
[Figure 3, “Basic Three-Room, Two-Node Campus Cluster Configuration With Multipathing,” on page 58](#) shows a three-room, two-node campus cluster.

FIGURE 3 Basic Three-Room, Two-Node Campus Cluster Configuration With Multipathing



In [Figure 4, “Minimum Three-Room, Two-Node Campus Cluster Configuration Without Multipathing,” on page 59](#), if at least two rooms are up and communicating, recovery is automatic. Only three-room or larger configurations can guarantee that the loss of any one room can be handled automatically.

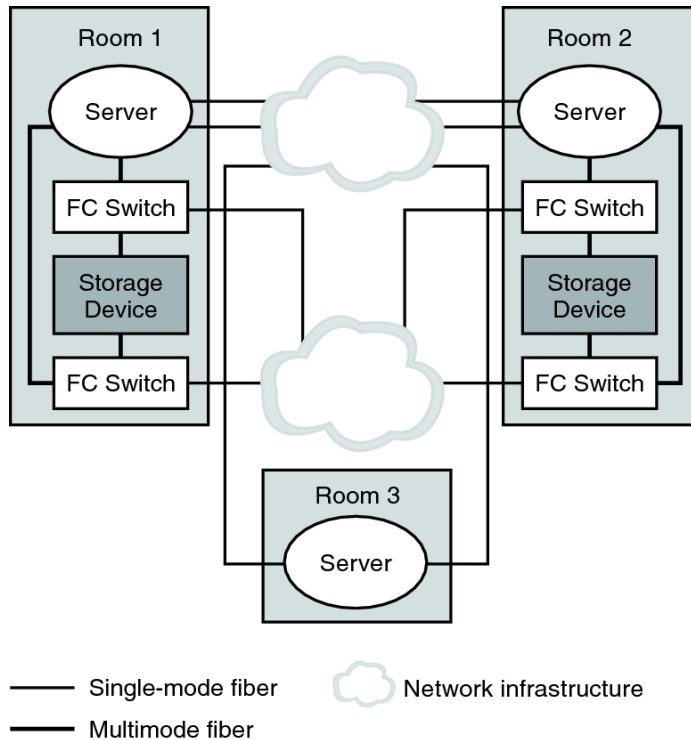
FIGURE 4 Minimum Three-Room, Two-Node Campus Cluster Configuration Without Multipathing



In [Figure 5, “Three-Room, Three-Node Campus Cluster Configuration,”](#) on page 60, one room contains one node and shared storage. A second room contains a cluster node only. The third room contains shared storage only. A LUN or disk of the storage device in the third room is configured as a quorum device.

This configuration provides the reliability of a three-room cluster with minimum hardware requirements. This campus cluster can survive the loss of any single room without requiring manual intervention.

FIGURE 5 Three-Room, Three-Node Campus Cluster Configuration



In [Figure 5, “Three-Room, Three-Node Campus Cluster Configuration,”](#) on page 60, a server acts as the quorum vote in the third room. This server does not necessarily support data services. Instead, it replaces a storage device as the quorum device.

Deciding How to Use Quorum Devices

When adding quorum devices to your campus cluster, your goal should be to balance the number of quorum votes in each room. No single room should have a much larger number of votes than the other rooms because loss of that room can bring the entire cluster down.

For campus clusters with more than three rooms and three nodes, quorum devices are optional. Whether you use quorum devices in such a cluster, and where you place them, depends on your assessment of the following:

- Your particular cluster topology
- The specific characteristics of the rooms involved
- Resiliency requirements for your cluster

As with two-room clusters, locate the quorum device in a room you determine is more likely to survive any failure scenario. Alternatively, you can locate the quorum device in a room that you *want* to form a cluster, in the event of a failure. Use your understanding of your particular cluster requirements to balance these two criteria.

Refer to the [Concepts for Oracle Solaris Cluster 4.4](#) for general information about quorum devices and how they affect clusters that experience failures. If you decide to use one or more quorum devices, consider the following recommended approach:

1. For each room, total the quorum votes (nodes) for that room.
2. Define a quorum device in the room that contains the lowest number of votes and that contains a fully connected shared storage device.

When your campus cluster contains more than two nodes, *do not* define a quorum device if each room contains the same number of nodes.

The following sections discuss quorum devices in various sizes of campus clusters.

- [“Quorum in Clusters With Four Rooms or More” on page 61](#)
- [“Quorum in Three-Room Configurations” on page 63](#)
- [“Quorum in Two-Room Configurations” on page 63](#)

Quorum in Clusters With Four Rooms or More

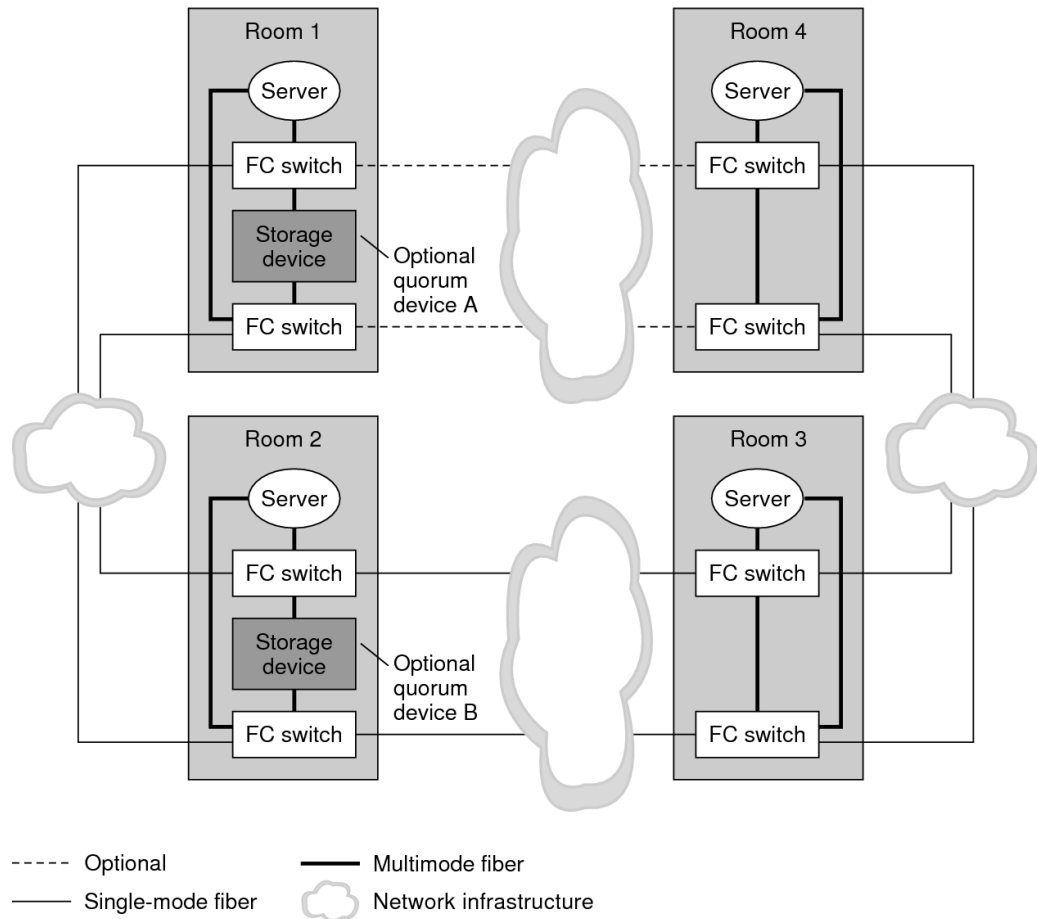
[Figure 6, “Four-Room, Four-Node Campus Cluster,” on page 62](#) illustrates a four-node campus cluster with fully connected storage. Each node is in a separate room. Two rooms also contain the shared storage devices, with data mirrored between them.

Note that the quorum devices are marked *optional* in the illustration. This cluster does not require a quorum device. With no quorum devices, the cluster can still survive the loss of any single room.

Consider the effect of adding *Quorum Device A*. Because the cluster contains four nodes, each with a single quorum vote, the quorum device receives three votes. Four votes (one node and the quorum device, or all four nodes) are required to form the cluster. This configuration is not optimal, because the loss of *Room 1* brings down the cluster. The cluster is not available after the loss of that single room.

If you then add *Quorum Device B*, both *Room 1* and *Room 2* have four votes. Six votes are required to form the cluster. This configuration is clearly better, as the cluster can survive the random loss of any single room.

FIGURE 6 Four-Room, Four-Node Campus Cluster



Note - In Figure 6, “Four-Room, Four-Node Campus Cluster,” on page 62, the cluster interconnect is not shown.

Consider the optional I/O connection between *Room 1* and *Room 4*. Although fully connected storage is preferable for reasons of redundancy and reliability, fully redundant connections might not always be possible in campus clusters. Geography might not accommodate a particular connection, or the project's budget might not cover the additional fiber.

In such a case, you can design a campus cluster with indirect access between some nodes and the storage. In Figure 7-4, if the optional I/O connection is omitted, *Node 4* must access the storage indirectly.

Quorum in Three-Room Configurations

In three-room, two-node campus clusters, you should use the third room for the quorum device ([Figure 3, “Basic Three-Room, Two-Node Campus Cluster Configuration With Multipathing,” on page 58](#)) or a server ([Figure 5, “Three-Room, Three-Node Campus Cluster Configuration,” on page 60](#)). Isolating the quorum device gives your cluster a better chance to maintain availability after the loss of one room. If at least one node and the quorum device remain operational, the cluster can continue to operate.

Quorum in Two-Room Configurations

In two-room configurations, the quorum device occupies the same room as one or more nodes. Place the quorum device in the room that is more likely to survive a failure scenario if all cluster transport and disk connectivity are lost between rooms. If *only* cluster transport is lost, the node that shares a room with the quorum device is not necessarily the node that reserves the quorum device first. For more information about quorum and quorum devices, see the [Concepts for Oracle Solaris Cluster 4.4](#).

Determining Campus Cluster Connection Technologies

This section lists example technologies for the private cluster interconnect and for the data paths and their various distance limits. In some cases, it is possible to extend these limits. For more information, ask your Oracle representative.

Cluster Interconnect Technologies

The following table lists example node-to-node link technologies and their limitations.

TABLE 7 Campus Cluster Interconnect Technologies and Distance Limits

Link Technology	Maximum Distance	Comments
100 Mbps Ethernet	100 meters per segment	Unshielded twisted pair (UTP)
1000 Mbps Ethernet	100 meters per segment	UTP
1000 Mbps Ethernet	260 meters per segment	62.5/125 micron multimode fiber (MMF)
1000 Mbps Ethernet	550 meters per segment	50/125 micron MMF
1000 Mbps Ethernet (FC)	10 kilometers at 1 Gbps	9/125 micron single-mode fiber (SMF)
DWDM	200 kilometers and up	
Other		Consult your Oracle representative

Always check your vendor documentation for technology-specific requirements and limitations.

Storage Area Network Technologies

The following table lists example link technologies for the cluster data paths and the distance limits for a single interswitch link (ISL).

TABLE 8 ISL Limits

Link Technology	Maximum Distance	Comments
FC short-wave gigabit interface converter (GBIC)	500 meters at 1 Gbps	50/125 micron MMF
FC long-wave GBIC	10 kilometers at 1 Gbps	9/125 micron SMF
FC short-wave small form-factor pluggable (SFP)	300 meters at 2 Gbps	62.5/125 micron MMF
FC short-wave SFP	500 meters at 2 Gbps	62.5/125 micron MMF
FC long-wave SFP	10 kilometers at 2 Gbps	9/125 micron SMF
DWDM	200 kilometers and up	
Other		Consult your Oracle representative

Installing and Configuring Interconnect, Storage, and Fibre Channel Hardware

Generally, using interconnect, storage, and Fibre Channel (FC) hardware does not differ markedly from standard cluster configurations.

The steps for installing Ethernet-based campus cluster interconnect hardware are the same as the steps for standard clusters. Refer to [“Installing Ethernet or InfiniBand Cluster Interconnect Hardware” on page 22](#) . When installing the media converters, consult the accompanying documentation, including requirements for fiber connections.

The guidelines for installing virtual local area networks interconnect networks are the same as the guidelines for standard clusters. See [“Configuring VLANs as Private Interconnect Networks” on page 24](#).

The steps for installing shared storage are the same as the steps for standard clusters. Refer to the documentation for your storage device for those steps.

Campus clusters require FC switches to mediate between multimode and single-mode fibers. The steps for configuring the settings on the FC switches are very similar to the steps for standard clusters.

If your switch supports flexibility in the buffer allocation mechanism, (for example the QLogic switch with donor ports), make certain you allocate a sufficient number of buffers to the ports that are dedicated to interswitch links (ISLs). If your switch has a fixed number of frame buffers (or buffer credits) per port, you do not have this flexibility.

Calculating Buffer Credits

The following rules determine the number of buffers that you might need:

- For 1 Gbps, calculate buffer credits as:

$$(length-in-km) \times (0.6)$$

Round the result up to the next whole number. For example, a 10 km connection requires 6 buffer credits, and a 7 km connection requires 5 buffer credits.

- For 2 Gbps, calculate buffer credits as:

$$(length-in-km) \times (1.2)$$

Round the result up to the next whole number. For example, a 10 km connection requires 12 buffer credits, while a 7 km connection requires 9 buffer credits.

For greater speeds or for more details, refer to your switch documentation for information about computing buffer credits.

Additional Campus Cluster Configuration Examples

While detailing all of the configurations that are possible in campus clustering is beyond the scope of this document, the following illustrations depict variations on the configurations that were previously shown.

- Three-room campus cluster with a multipathing solution implemented ([Figure 7, “Three-Room Campus Cluster With a Multipathing Solution Implemented,”](#) on page 67)
- Two-room campus cluster with a multipathing solution implemented ([Figure 8, “Two-Room Campus Cluster With a Multipathing Solution Implemented,”](#) on page 68)
- Two-room campus cluster without a multipathing solution implemented ([Figure 9, “Two-Room Campus Cluster Without a Multipathing Solution Implemented,”](#) on page 69)

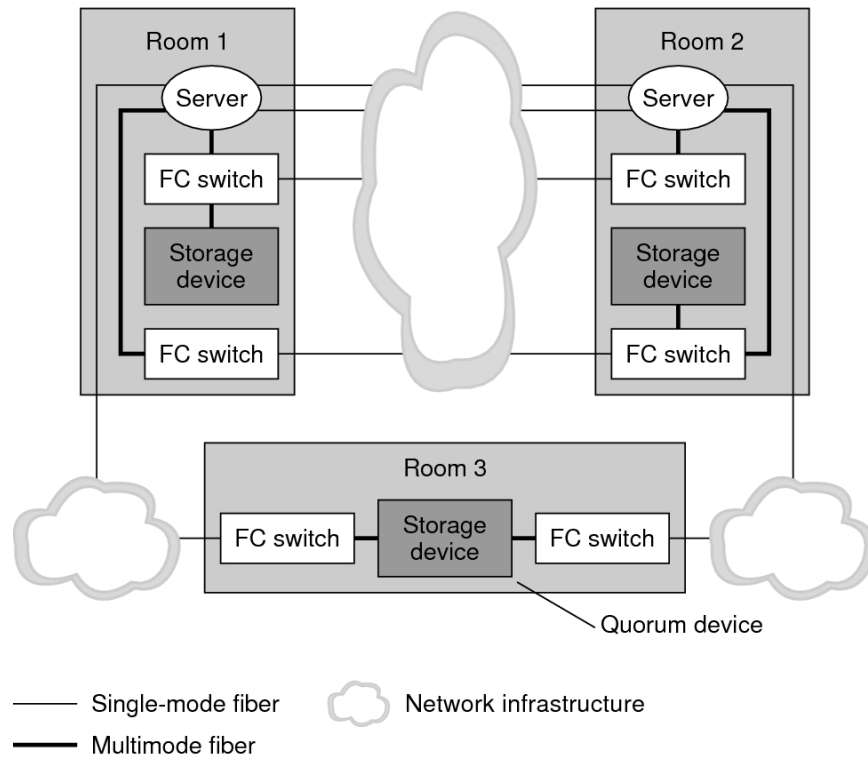
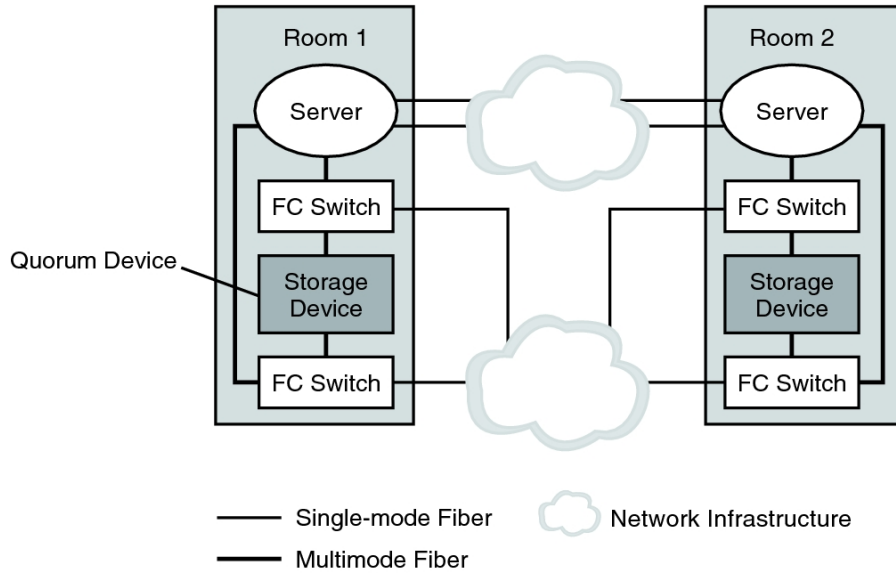
FIGURE 7 Three-Room Campus Cluster With a Multipathing Solution Implemented

Figure 8, “Two-Room Campus Cluster With a Multipathing Solution Implemented,” on page 68 shows a two-room campus cluster that uses partner pairs of storage devices and four FC switches, with a multipathing solution implemented. The four switches are added to the cluster for greater redundancy and potentially better I/O throughput. Other possible configurations that you could implement include using Oracle’s Sun StorEdge T3 partner groups or Oracle’s Sun StorEdge 9910/9960 arrays with Oracle Solaris I/O multipathing software installed.

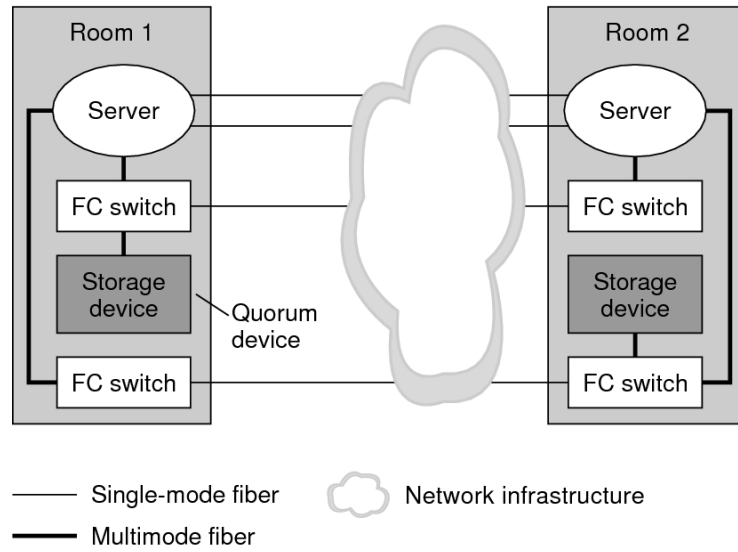
For information about Solaris I/O multipathing software for the Oracle Solaris 11 OS, see the [Managing SAN Devices and I/O Multipathing in Oracle Solaris 11.4](#).

FIGURE 8 Two-Room Campus Cluster With a Multipathing Solution Implemented



The configuration in the following figure could be implemented by using Oracle's Sun StorEdge T3 or T3+ arrays in single-controller configurations, rather than partner groups.

FIGURE 9 Two-Room Campus Cluster Without a Multipathing Solution Implemented



Verifying Oracle Solaris Cluster Hardware Redundancy

This chapter describes the tests for verifying and demonstrating the high availability (HA) of your Oracle Solaris Cluster configuration. The tests in this chapter assume that you installed Oracle Solaris Cluster hardware, the Oracle Solaris Operating System, and Oracle Solaris Cluster software. All nodes should be booted as cluster members.

This chapter contains the following procedures:

- [“How to Test Device Group Redundancy Using Resource Group Failover” on page 72](#)
- [“How to Test Cluster Interconnects” on page 73](#)
- [“How to Test Public Network Redundancy” on page 74](#)

If your cluster passes these tests, your hardware has adequate redundancy. This redundancy means that your nodes, cluster transport cables, and IPMP groups are not single points of failure.

To perform the tests in [“How to Test Device Group Redundancy Using Resource Group Failover” on page 72](#) and [“How to Test Cluster Interconnects” on page 73](#), you must first identify the device groups that each node masters. Perform these tests on all cluster pairs that share a disk device group. Each pair has a primary node and a secondary node for a particular device group.

Use the following command to determine the initial primary and secondary: `cldevicegroup` status with the `-n` option.

For conceptual information about primary nodes, secondary nodes, failover, device groups, or cluster hardware, see the [Concepts for Oracle Solaris Cluster 4.4](#).

Testing Node Redundancy

This section provides the procedure for testing node redundancy and high availability of device groups. Perform the following procedure to confirm that the secondary node takes over the device group that is mastered by the primary node when the primary node fails.

▼ How to Test Device Group Redundancy Using Resource Group Failover

Before You Begin To perform this procedure, become an administrator that provides `solaris.cluster.modify` authorization.

1. Create an HASStoragePlus resource group with which to test.

Use the following command:

```
# clresourcegroup create testgroup
# clresourcetype register SUNW.HASStoragePlus
# clresource create -t HASStoragePlus -g testgroup \
-p GlobalDevicePaths=/dev/md/red/dsk/d0 \
-p Affinityon=true testresource
```

`clresourcetype register` If the HASStoragePlus resource type is not already registered, register it.

`/dev/md/red/dsk/d0` Replace this path with your device path.

2. Identify the node that masters the testgroup.

```
# clresourcegroup status testgroup
```

3. Power off the primary node for the testgroup.

Cluster interconnect error messages appear on the consoles of the existing nodes.

4. On another node, verify that the secondary node took ownership of the resource group that is mastered by the primary node.

Use the following command to check the output for the resource group ownership:

```
# clresourcegroup status testgroup
```

5. Power on the initial primary node. Boot the node into cluster mode.

Wait for the system to boot. The system automatically starts the membership monitor software. The node then rejoins the cluster.

6. **From the initial primary node, return ownership of the resource group to the initial primary node.**

```
# clresourcegroup switch -n nodename testgroup
```

In these commands, *nodename* is the name of the primary node.

7. **Verify that the initial primary node has ownership of the resource group.**
Use the following command to look for the output that shows the device group ownership.

```
# clresourcegroup status testgroup
```

Testing Cluster Interconnect Redundancy

This section provides the procedure for testing cluster interconnect redundancy.

▼ How to Test Cluster Interconnects

Before You Begin To perform this procedure, become an administrator that provides `solaris.cluster.read` and `solaris.cluster.modify` authorization.

1. **Disconnect one of the cluster transport cables from a node in the cluster.**

Messages similar to the following appear on the consoles of each node and are logged in the `/var/adm/messages` file.

```
Nov  4 08:27:21 node1 genunix: WARNING: ce1: fault detected external to device; service
degraded
Nov  4 08:27:21 node1 genunix: WARNING: ce1: xcvr addr:0x01 - link down
Nov  4 08:27:31 node1 genunix: NOTICE: clcomm: Path node1:ce1 - node1:ce0 being cleaned
up
Nov  4 08:27:31 node1 genunix: NOTICE: clcomm: Path node1:ce1 - node1:ce0 being drained
Nov  4 08:27:31 node1 genunix: NOTICE: clcomm: Path node1:ce1 - node1:ce0 being
constructed
Nov  4 08:28:31 node1 genunix: NOTICE: clcomm: Path node1:ce1 - node1:ce0 errors during
initiation
Nov  4 08:28:31 node1 genunix: WARNING: Path node1:ce1 - node1:ce0 initiation
encountered errors, errno = 62.
Remote node may be down or unreachable through this path.
```

2. Verify that Oracle Solaris Cluster has registered that the interconnect is down.

Use the following command to verify that the interconnect path displays as `Faulted`.

```
# clinterconnect status
```

3. Reconnect the cluster transport cable.

Messages similar to the following appear on the consoles of each node and are logged in the `/var/adm/messages` file.

```
Nov  4 08:30:26 node1 genunix: NOTICE: ce1: fault cleared external to device; service
  available
Nov  4 08:30:26 node1 genunix: NOTICE: ce1: xcvr addr:0x01 - link up 1000 Mbps full
  duplex
Nov  4 08:30:26 node1 genunix: NOTICE: clcomm: Path node1:ce1 - node1:ce0 being
  initiated
Nov  4 08:30:26 node1 genunix: NOTICE: clcomm: Path node1:ce1 - node1:ce0 online
```

4. Verify that Oracle Solaris Cluster has registered that the interconnect is up.

Use the following command to verify that the interconnect path displays as `Online`.

```
# clinterconnect status
```

5. Repeat [Step 1](#) through [Step 4](#) on each cluster transport cable in the node.

6. Repeat [Step 1](#) through [Step 5](#) on each node in the cluster.

Testing Public Network Redundancy

This section provides the procedure for testing public network redundancy.

▼ How to Test Public Network Redundancy

If you perform this test, you can verify that IP addresses failover from one adapter to another adapter within the same IPMP group.

Ensure that the `/etc/netmasks` file has IP-address subnet and netmask entries for all logical hostnames. If necessary, edit the `/etc/netmasks` file to add any missing entries.

Before You Begin To perform this procedure, become an administrator that provides `solaris.cluster.read` authorization.

1. **Create a logical hostname resource group which is the failover hostname to use the IPMP groups on the system.**

Use the following command:

```
# clresourcegroup create lhtestgroup
# clreslogicalhostname create -g lhtestgroup logicalhostname
# clresourcegroup online lhtestgroup
```

logicalhostname The IP address that is hosted on the device on which an IPMP group is configured.

2. **Determine the adapter on which the *logicalhostname* exists.**

```
# ipadm show-if -o all
```

3. **Disconnect one public network cable from the adapter you identified in [Step 2](#).**
4. **If there are no more adapters in the group, skip to [Step 7](#).**
5. **If there is another adapter in the group, verify that the logical hostname failed over to that adapter.**

```
# ipadm show-if -o all
```

6. **Continue to disconnect adapters in the group, until you have disconnected the last adapter.**

The resource group (*lhtestgroup*) should fail over to the secondary node.

7. **Verify that the resource group failed over to the secondary node.**

Use the following command:

```
# clnode status lhtestgroup
```

8. **Reconnect all adapters in the group.**
9. **From the initial primary node, return ownership of the resource group to the initial primary node.**

```
# clresourcegroup switch -n nodename lhtestgroup
```

In these commands, *nodename* is the name of the original primary node.

10. **Verify that the resource group is running on the original primary node.**

Use the following command:

```
# clnode status lhstestgroup
```

Index

A

adapters *See* public network adapters *See* transport adapters
adding
 public network adapters, 39
 transport adapters, 30
 transport cables, 30
 transport junctions, 30

B

boot, 14

C

cables *See* transport cables
campus clustering, 53
 and fibre channel, 65
 configuration examples, 57, 57
 data replication requirements, 55
 differences from standard, 53
 extended examples, 66
 guidelines, 56
 hardware installation and configuration, 65
 interconnect technologies, 63
 link technologies, 63
 network requirements, 54
 node compared to room, 56
 private network technologies, 63
 public network technologies, 64
 quorum guidelines, 60
 quorum requirements, 55
 requirements, 53

 room, definition of, 56
 three-room examples, 67
 two-room, example without multipathing, 68
 two-room, multipathing example, 67
 VLANs, 26
 volume management, 55
CD-ROMs, 16
cldevice
 populate, 14
clnode
 evacuate, 14
cluster
 interconnect, 13
concurrent access, 13
configuration
 parallel database, 13
configuration examples (campus clustering)
 three-room, 57, 57

D

data replication
 requirements, 55
dual-port HBAs, 45
 Solaris Volume Manager, 47
 Solaris Volume Manager for Sun Cluster, 48
 supported configurations, 47
dynamic reconfiguration
 kernel cage recovery, 49
 preparing the cluster for kernel cage dynamic reconfiguration, 49
 recovering from an interrupted kernel cage dynamic reconfiguration, 50

- replacing disks, 15
- E**
 - Ethernet switches
 - in the interconnect, 20
 - Ethernet transport cables and junctions, 23
 - example configurations (campus clustering)
 - three-room, 57
- F**
 - fibre channel and campus clustering, 65
- G**
 - GUI
 - tasks you can perform
 - add an interconnect component, 30
- H**
 - hardware installation
 - for campus clusters, 65
 - overall installation process, 13
 - overview, 11
 - hardware RAID, 41
 - hardware redundancy
 - verifying, 71
 - hardware restrictions, 17
 - high availability
 - testing, 71
 - testing device group availability, 72
 - testing interconnect availability, 73
 - testing IP multipathing availability, 74
 - host adapters
 - dual-port configurations, 45
- I**
 - InfiniBand
 - Socket Direct Protocol, 22
 - InfiniBand requirements and restrictions, 21
 - InfiniBand transport cables and junctions, 23
 - installing
 - cluster hardware, 11
 - Ethernet transport cables, 23
 - Ethernet transport junctions, 23
 - InfiniBand transport cables, 23
 - InfiniBand transport junctions, 23
 - Oracle Solaris and cluster software, 14
 - public network hardware, 38
 - integrated mirroring, 41
 - interconnect
 - configuration for campus clustering, 65
 - jumbo frames requirements, 30
 - speed requirements, 20
 - technologies for campus clustering, 63
 - testing redundancy, 73
 - internal hardware disk mirroring, 41
 - IP multipathing
 - testing redundancy, 74
- J**
 - jumbo frames
 - interconnect requirements, 30
 - public network requirements, 20, 37
 - Scalable Data Services, 37
- K**
 - kernel cage dynamic reconfiguration
 - preparing the cluster, 49
 - recovering from an interruption, 50
 - recovery, 49
- L**
 - link technologies
 - campus clustering, 63
 - local disks, 16

M

- mirroring internal disks, 41
- multihost disks, 16
- multipathing
 - example three-room campus cluster, 67
 - example two-room campus cluster, 67

N

- NAFO groups
 - adding adapters, 40
 - redundancy testing, 74
- network
 - private, 13
- Network Adapter Failover groups *See* NAFO groups
- network requirements for campus clusters, 54
- node redundancy
 - testing, 72
- nodes
 - must be of the same architecture, 17

O

- Oracle Real Application Clusters, 54

P

- parallel database configurations, 13
- powering off, 15
- powering on, 15
- private network, 13
- private network technologies
 - campus clustering, 63
- public network
 - hardware installation, 38
 - jumbo frames requirements, 20, 37
- public network adapters
 - adding, 39
 - removing, 39, 40
 - replacing, 39
- public network technologies

- campus clustering, 64

Q

- quorum devices
 - campus cluster guidelines, 60
 - campus cluster requirements, 55

R

- raidctl command, 41
- redundancy
 - testing interconnect redundancy, 73
 - testing IP multipathing redundancy, 74
 - testing node redundancy, 72
- removable media, 16
- removing
 - public network adapters, 39, 40
 - transport adapters, 33
 - transport cables, 33
 - transport junctions, 33
- replacing
 - public network adapters, 39
 - transport adapters, 31
 - transport cables, 31
 - transport junctions, 31
- requirements
 - interconnect speed, 20
- restrictions *See* hardware
- room compared to node (campus clustering), 56
- room, definition of (campus clustering), 56

S

- SAN
 - general cluster requirements, 17
 - requirements in campus clusters, 55
- SDP *See* Socket Direct Protocol
- shutdown, 14
- shutdown protocol
 - clustered environment, 15
 - nonclustered environment, 15

- Socket Direct Protocol, 22
- software installation, 14
- Solaris Volume Manager
 - dual-port HBAs, 47
- Solaris Volume Manager for Sun Cluster
 - dual-port HBAs, 48
- SR-IOV devices, 27
- standard clusters
 - differences from campus clusters, 53
- supported configurations, 12
- switches, 23 *See* transport junctions
 - See also* transport junctions

T

- tapes, 16
- testing
 - high availability, 71
 - interconnect availability, 73
 - interconnect redundancy, 73
 - IP multipathing availability, 74
 - IP multipathing redundancy, 74
 - NAFO group redundancy, 74
 - node availability, 72
 - node redundancy, 72
- transport adapter firmware
 - upgrading, 35
- transport adapters
 - adding, 30
 - removing, 33
 - replacing, 31
- transport cables
 - adding, 30
 - Ethernet, installing, 23
 - InfiniBand, installing, 23
 - removing, 33
 - replacing, 31
- transport junctions
 - adding, 30
 - Ethernet, installing, 23
 - InfiniBand, installing, 23
 - removing, 33
 - replacing, 31

U

- upgrading
 - transport adapter firmware, 35

V

- verifying
 - hardware redundancy, 71
- virtual local area networks *See* VLANs
- VLANs
 - campus clusters, 26
 - configuring SR-IOV devices, 27
 - guidelines, 24
 - volume management with campus clustering, 55