

Veritas Storage Foundation™ for Oracle® RAC 6.0.1 Administrator's Guide - HP-UX

Veritas Storage Foundation™ for Oracle RAC Administrator's Guide

The software described in this book is furnished under a license agreement and may be used only in accordance with the terms of the agreement.

Product version: 6.0.1

Document version: 6.0.1 Rev 1

Legal Notice

Copyright © 2012 Symantec Corporation. All rights reserved.

Symantec, the Symantec logo, Veritas, Veritas Storage Foundation, CommandCentral, NetBackup, Enterprise Vault, and LiveUpdate are trademarks or registered trademarks of Symantec corporation or its affiliates in the U.S. and other countries. Other names may be trademarks of their respective owners.

The product described in this document is distributed under licenses restricting its use, copying, distribution, and decompilation/reverse engineering. No part of this document may be reproduced in any form by any means without prior written authorization of Symantec Corporation and its licensors, if any.

THE DOCUMENTATION IS PROVIDED "AS IS" AND ALL EXPRESS OR IMPLIED CONDITIONS, REPRESENTATIONS AND WARRANTIES, INCLUDING ANY IMPLIED WARRANTY OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE OR NON-INFRINGEMENT, ARE DISCLAIMED, EXCEPT TO THE EXTENT THAT SUCH DISCLAIMERS ARE HELD TO BE LEGALLY INVALID. SYMANTEC CORPORATION SHALL NOT BE LIABLE FOR INCIDENTAL OR CONSEQUENTIAL DAMAGES IN CONNECTION WITH THE FURNISHING, PERFORMANCE, OR USE OF THIS DOCUMENTATION. THE INFORMATION CONTAINED IN THIS DOCUMENTATION IS SUBJECT TO CHANGE WITHOUT NOTICE.

The Licensed Software and Documentation are deemed to be commercial computer software as defined in FAR 12.212 and subject to restricted rights as defined in FAR Section 52.227-19 "Commercial Computer Software - Restricted Rights" and DFARS 227.7202, "Rights in Commercial Computer Software or Commercial Computer Software Documentation", as applicable, and any successor regulations. Any use, modification, reproduction release, performance, display or disclosure of the Licensed Software and Documentation by the U.S. Government shall be solely in accordance with the terms of this Agreement.

Symantec Corporation
350 Ellis Street
Mountain View, CA 94043
<http://www.symantec.com>

Technical Support

Symantec Technical Support maintains support centers globally. Technical Support's primary role is to respond to specific queries about product features and functionality. The Technical Support group also creates content for our online Knowledge Base. The Technical Support group works collaboratively with the other functional areas within Symantec to answer your questions in a timely fashion. For example, the Technical Support group works with Product Engineering and Symantec Security Response to provide alerting services and virus definition updates.

Symantec's support offerings include the following:

- A range of support options that give you the flexibility to select the right amount of service for any size organization
- Telephone and/or Web-based support that provides rapid response and up-to-the-minute information
- Upgrade assurance that delivers software upgrades
- Global support purchased on a regional business hours or 24 hours a day, 7 days a week basis
- Premium service offerings that include Account Management Services

For information about Symantec's support offerings, you can visit our Web site at the following URL:

www.symantec.com/business/support/index.jsp

All support services will be delivered in accordance with your support agreement and the then-current enterprise technical support policy.

Contacting Technical Support

Customers with a current support agreement may access Technical Support information at the following URL:

www.symantec.com/business/support/contact_techsupp_static.jsp

Before contacting Technical Support, make sure you have satisfied the system requirements that are listed in your product documentation. Also, you should be at the computer on which the problem occurred, in case it is necessary to replicate the problem.

When you contact Technical Support, please have the following information available:

- Product release level

- Hardware information
- Available memory, disk space, and NIC information
- Operating system
- Version and patch level
- Network topology
- Router, gateway, and IP address information
- Problem description:
 - Error messages and log files
 - Troubleshooting that was performed before contacting Symantec
 - Recent software configuration changes and network changes

Licensing and registration

If your Symantec product requires registration or a license key, access our technical support Web page at the following URL:

www.symantec.com/business/support/

Customer service

Customer service information is available at the following URL:

www.symantec.com/business/support/

Customer Service is available to assist with non-technical questions, such as the following types of issues:

- Questions regarding product licensing or serialization
- Product registration updates, such as address or name changes
- General product information (features, language availability, local dealers)
- Latest information about product updates and upgrades
- Information about upgrade assurance and support contracts
- Information about the Symantec Buying Programs
- Advice about Symantec's technical support options
- Nontechnical presales questions
- Issues that are related to CD-ROMs or manuals

Support agreement resources

If you want to contact Symantec regarding an existing support agreement, please contact the support agreement administration team for your region as follows:

Asia-Pacific and Japan	customercare_apac@symantec.com
Europe, Middle-East, and Africa	semea@symantec.com
North America and Latin America	supportsolutions@symantec.com

Documentation

Product guides are available on the media in PDF format. Make sure that you are using the current version of the documentation. The document version appears on page 2 of each guide. The latest product documentation is available on the Symantec Web site.

<https://sort.symantec.com/documents>

Your feedback on product documentation is important to us. Send suggestions for improvements and reports on errors or omissions. Include the title and document version (located on the second page), and chapter and section titles of the text on which you are reporting. Send feedback to:

doc_feedback@symantec.com

For information regarding the latest HOWTO articles, documentation updates, or to ask a question regarding product documentation, visit the Storage and Clustering Documentation forum on Symantec Connect.

<https://www-secure.symantec.com/connect/storage-management/forums/storage-and-clustering-documentation>

About Symantec Connect

Symantec Connect is the peer-to-peer technical community site for Symantec's enterprise customers. Participants can connect and share information with other product users, including creating forum posts, articles, videos, downloads, blogs and suggesting ideas, as well as interact with Symantec product teams and Technical Support. Content is rated by the community, and members receive reward points for their contributions.

<http://www.symantec.com/connect/storage-management>

Applying Oracle patches on SF Oracle RAC nodes	102
Installing Veritas Volume Manager, Veritas File System, or ODM patches on SF Oracle RAC nodes	103
Applying operating system updates on SF Oracle RAC nodes	103
Determining the LMX traffic for each database in an SF Oracle RAC cluster	104
Adding storage to an SF Oracle RAC cluster	106
Recovering from storage failure	107
Backing up and restoring Oracle database using Symantec NetBackup	107
Enhancing the performance of SF Oracle RAC clusters	108
Creating snapshots for offhost processing	109
Managing database storage efficiently using SmartTier	109
Optimizing database storage using Thin Provisioning and SmartMove	109
Scheduling periodic health checks for your SF Oracle RAC cluster	109
Verifying the nodes in an SF Oracle RAC cluster	110
Administering VCS	111
Viewing available Veritas device drivers	112
Loading Veritas drivers into memory	112
Starting and stopping VCS	112
Environment variables to start and stop VCS modules	113
Adding and removing LLT links	115
Configuring aggregated interfaces under LLT	119
Displaying the cluster details and LLT version for LLT links	121
Configuring destination-based load balancing for LLT	122
Enabling and disabling intelligent resource monitoring for agents manually	122
Administering the AMF kernel driver	124
Administering I/O fencing	125
About administering I/O fencing	126
About the vxfcntl utility	127
About the vxfcntladm utility	134
About the vxfcntlpre utility	139
About the vxfcntlswap utility	142
Enabling or disabling the preferred fencing policy	154
Administering the CP server	156
Refreshing registration keys on the coordination points for server-based fencing	156
Replacing coordination points for server-based fencing in an online cluster	158

	Migrating from non-secure to secure setup for CP server and SF Oracle RAC cluster communication	161
	Administering CFS	162
	Resizing CFS file systems	163
	Verifying the status of CFS file system nodes and their mount points	163
	Administering CVM	164
	Listing all the CVM shared disks	164
	Establishing CVM cluster membership manually	165
	Changing the CVM master manually	165
	Importing a shared disk group manually	168
	Deporting a shared disk group manually	169
	Starting shared volumes manually	169
	Verifying if CVM is running in an SF Oracle RAC cluster	169
	Verifying CVM membership state	170
	Verifying the state of CVM shared disk groups	170
	Verifying the activation mode	170
	Administering SF Oracle RAC global clusters	171
	About setting up a disaster recovery fire drill	171
	About configuring the fire drill service group using the Fire Drill Setup wizard	172
	Verifying a successful fire drill	174
	Scheduling a fire drill	174
	Sample fire drill service group configuration	174
Section 2	Performance and troubleshooting	177
Chapter 3	Troubleshooting SF Oracle RAC	179
	About troubleshooting SF Oracle RAC	179
	Gathering information from an SF Oracle RAC cluster for support analysis	180
	SF Oracle RAC log files	182
	About SF Oracle RAC kernel and driver messages	185
	VCS message logging	186
	What to do if you see a licensing reminder	191
	Restarting the installer after a failed connection	192
	Installer cannot create UUID for the cluster	192
	Troubleshooting SF Oracle RAC pre-installation check failures	193
	Troubleshooting LLT health check warning messages	195
	Troubleshooting LMX health check warning messages in SF Oracle RAC clusters	198
	Troubleshooting I/O fencing	199

	Oracle log files show shutdown called even when not shutdown manually	222
	Resolving ASYNCH_IO errors in an SF Oracle RAC cluster	222
	Oracle's clusterware processes fail to start	223
	Oracle Clusterware fails after restart	223
	Troubleshooting the Virtual IP (VIP) configuration in an SF Oracle RAC cluster	224
	Troubleshooting Oracle Clusterware health check warning messages in SF Oracle RAC clusters	225
	Troubleshooting ODM in SF Oracle RAC clusters	226
	File System configured incorrectly for ODM shuts down Oracle	226
Chapter 4	Prevention and recovery strategies	229
	Verification of GAB ports in SF Oracle RAC cluster	229
	Examining GAB seed membership	230
	Manual GAB membership seeding	231
	Evaluating VCS I/O fencing ports	232
	Verifying normal functioning of VCS I/O fencing	233
	Managing SCSI-3 PR keys in SF Oracle RAC cluster	233
	Evaluating the number of SCSI-3 PR keys on a coordinator LUN, if there are multiple paths to the LUN from the hosts	234
	Detecting accidental SCSI-3 PR key removal from coordinator LUNs	234
	Identifying a faulty coordinator LUN	235
Chapter 5	Tunable parameters	237
	About SF Oracle RAC tunable parameters	237
	About GAB tunable parameters	238
	About GAB load-time or static tunable parameters	238
	About GAB run-time or dynamic tunable parameters	240
	About LLT tunable parameters	245
	About LLT timer tunable parameters	246
	About LLT flow control tunable parameters	250
	Setting LLT timer tunable parameters	252
	About LMX tunable parameters	253
	LMX tunable parameters	253
	About VXFEN tunable parameters	256
	Configuring the VXFEN module parameters	257
	Tuning guidelines for campus clusters	258

Section 3	Reference	259
Appendix A	List of SF Oracle RAC health checks	261
	LLT health checks	261
	LMX health checks in SF Oracle RAC clusters	265
	I/O fencing health checks	266
	PrivNIC health checks in SF Oracle RAC clusters	268
	Oracle Clusterware health checks in SF Oracle RAC clusters	269
	CVM, CFS, and ODM health checks in SF Oracle RAC clusters	270
Appendix B	Error messages	271
	About error messages	271
	LMX error messages in SF Oracle RAC	271
	LMX critical error messages in SF Oracle RAC	272
	LMX non-critical error messages in SF Oracle RAC	273
	VxVM error messages	274
	VXFEN driver error messages	274
	VXFEN driver informational message	275
	Node ejection informational messages	275
Glossary	277
Index	281

SF Oracle RAC concepts and administration

- [Chapter 1. Overview of Veritas Storage Foundation for Oracle RAC](#)
- [Chapter 2. Administering SF Oracle RAC and its components](#)

Overview of Veritas Storage Foundation for Oracle RAC

This chapter includes the following topics:

- [About Veritas Storage Foundation for Oracle RAC](#)
- [How SF Oracle RAC works \(high-level perspective\)](#)
- [Component products and processes of SF Oracle RAC](#)
- [Periodic health evaluation of SF Oracle RAC clusters](#)
- [About Virtual Business Services](#)
- [About Veritas Operations Manager](#)
- [About Symantec Operations Readiness Tools](#)

About Veritas Storage Foundation for Oracle RAC

Veritas Storage Foundation™ for Oracle® RAC (SF Oracle RAC) leverages proprietary storage management and high availability technologies to enable robust, manageable, and scalable deployment of Oracle RAC on UNIX platforms. The solution uses Veritas Cluster File System technology that provides the dual advantage of easy file system management as well as the use of familiar operating system tools and utilities in managing databases.

The solution stack comprises the Veritas Cluster Server (VCS), Veritas Cluster Volume Manager (CVM), Veritas Oracle Real Application Cluster Support (VRTSdbac), Veritas Oracle Disk Manager (VRTSodm), Veritas Cluster File System (CFS), and Veritas Storage Foundation, which includes the base Veritas Volume Manager (VxVM) and Veritas File System (VxFS).

Benefits of SF Oracle RAC

SF Oracle RAC provides the following benefits:

- Support for file system-based management. SF Oracle RAC provides a generic clustered file system technology for storing and managing Oracle data files as well as other application data.
- Support for high-availability of cluster interconnects.
For Oracle RAC 10g Release 2:
The combination of LMX/LLT protocols and the PrivNIC/MultiPrivNIC agents provides maximum bandwidth as well as high availability of the cluster interconnects, including switch redundancy.
For Oracle RAC 11g Release 1/Oracle RAC 11g Release 2:
The PrivNIC/MultiPrivNIC agents provide maximum bandwidth as well as high availability of the cluster interconnects, including switch redundancy.
See the following Technote regarding co-existence of PrivNIC/MultiPrivNIC agents with Oracle RAC 11.2.0.2 and later versions:
<http://www.symantec.com/business/support/index?page=content&id=TECH145261>
- Use of Cluster File System and Cluster Volume Manager for placement of Oracle Cluster Registry (OCR) and voting disks. These technologies provide robust shared block interfaces (for all supported Oracle RAC versions) and raw interfaces (for Oracle RAC 10g Release 2) for placement of OCR and voting disks.
- Support for a standardized approach toward application and database management. Administrators can apply their expertise of Veritas technologies toward administering SF Oracle RAC.
- Increased availability and performance using Veritas Dynamic Multi-Pathing (DMP). DMP provides wide storage array support for protection from failures and performance bottlenecks in the Host Bus Adapters (HBA), Storage Area Network (SAN) switches, and storage arrays.
- Easy administration and monitoring of multiple SF Oracle RAC clusters using Veritas Operations Manager.
- VCS OEM plug-in provides a way to monitor SF Oracle RAC resources from the OEM console.
For more information, see the *Veritas Storage Foundation: Storage and Availability Management for Oracle Databases* guide.
- Improved file system access times using Oracle Disk Manager (ODM).
- Ability to configure Oracle Automatic Storage Management (ASM) disk groups over CVM volumes to take advantage of Veritas Dynamic Multi-Pathing (DMP).

- Enhanced scalability and availability with access to multiple Oracle RAC instances per database in a cluster.
- Support for backup and recovery solutions using volume-level and file system-level snapshot technologies, Storage Checkpoints, and Database Storage Checkpoints.
 For more information, see the *Veritas Storage Foundation: Storage and Availability Management for Oracle Databases* guide.
- Support for space optimization using periodic deduplication in a file system to eliminate duplicate data without any continuous cost.
 For more information, see the Veritas Storage Foundation Administrator's documentation.
- Ability to fail over applications with minimum downtime using Veritas Cluster Server (VCS) and Veritas Cluster File System (CFS).
- Prevention of data corruption in split-brain scenarios with robust SCSI-3 Persistent Group Reservation (PGR) based I/O fencing or Coordination Point Server-based I/O fencing. The preferred fencing feature also enables you to specify how the fencing driver determines the surviving subcluster.
- Support for sharing application data, in addition to Oracle database files, across nodes.
- Support for policy-managed databases in Oracle RAC 11g Release 2.
- Fast disaster recovery with minimal downtime and interruption to users. Users can transition from a local high availability site to a wide-area disaster recovery environment with primary and secondary sites. If a site fails, clients that are attached to the failed site can reconnect to a surviving site and resume access to the shared database.
- Verification of disaster recovery configuration using fire drill technology without affecting production systems.
- Support for a wide range of hardware replication technologies as well as block-level replication using VVR.
- Support for campus clusters with the following capabilities:
 - Consistent detach with Site Awareness
 - Site aware reads with VxVM mirroring
 - Monitoring of Oracle resources
 - Protection against split-brain scenarios

How SF Oracle RAC works (high-level perspective)

Oracle Real Application Clusters (RAC) is a parallel database environment that takes advantage of the processing power of multiple computers. Oracle stores data logically in the form of tablespaces and physically in the form of data files. The Oracle instance is a set of processes and shared memory that provide access to the physical database. Specifically, the instance involves server processes acting on behalf of clients to read data into shared memory and make modifications to it, and background processes that interact with each other and with the operating system to manage memory structure and do general housekeeping.

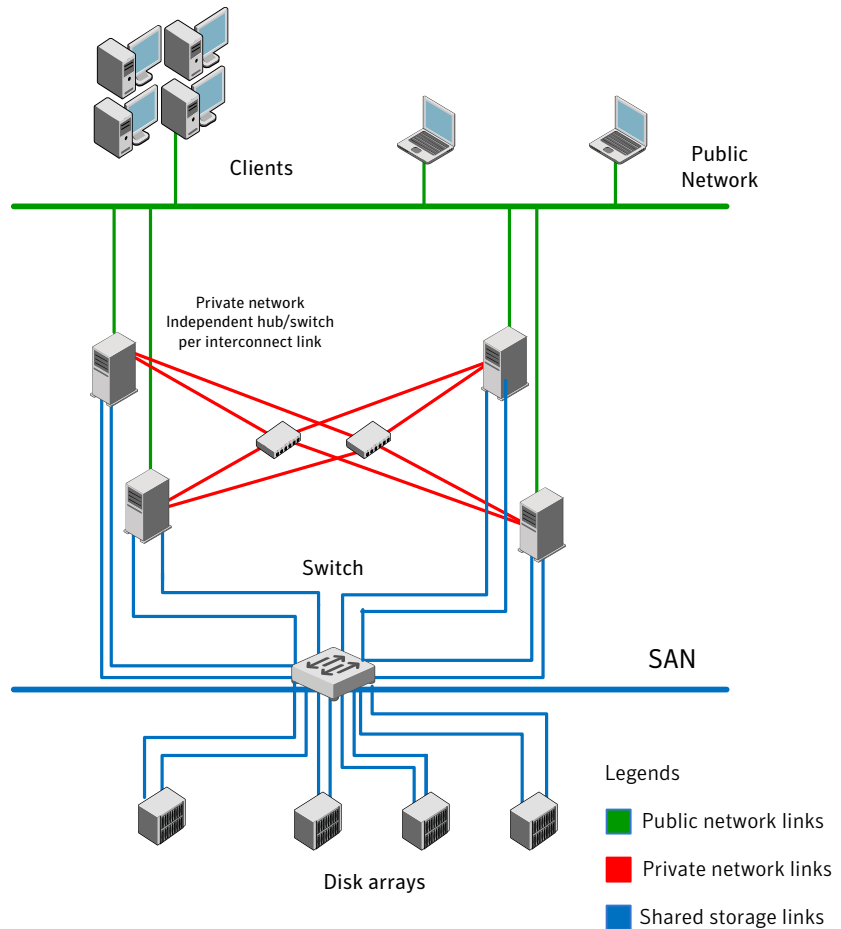
SF Oracle RAC provides the necessary infrastructure for running Oracle RAC and coordinates access to the shared data for each node to provide consistency and integrity. Each node adds its processing power to the cluster as a whole and can increase overall throughput or performance.

At a conceptual level, SF Oracle RAC is a cluster that manages applications (Oracle instances), networking, and storage components using resources contained in service groups. SF Oracle RAC clusters have the following properties:

- Each node runs its own operating system.
- A cluster interconnect enables cluster communications.
- A public network connects each node to a LAN for client access.
- Shared storage is accessible by each node that needs to run the application.

[Figure 1-1](#) displays the basic layout and individual components required for a SF Oracle RAC installation.

Figure 1-1 SF Oracle RAC basic layout and components



The basic layout has the following characteristics:

- Multiple client applications that access nodes in the cluster over a public network.
- Nodes that are connected by at least two private network links (also called cluster interconnects) using 100BaseT or gigabit Ethernet controllers on each system.
If the private links are on a single switch, isolate them using VLAN.
- Nodes that are connected to iSCSI or Fibre Channel shared storage devices over SAN.
All shared storage must support SCSI-3 PR.

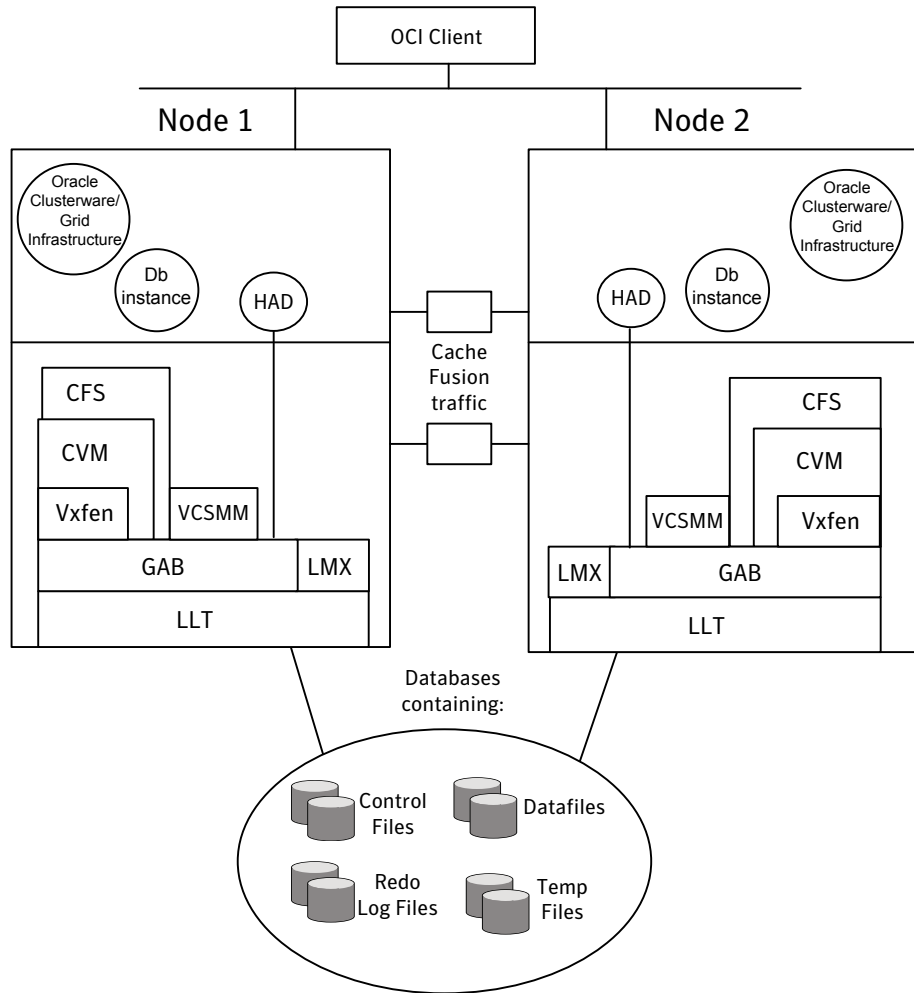
- Nodes must be connected with private network links using similar network devices.
- The Oracle Cluster Registry, vote disks, and data files configured on the shared storage that is available to each node.
For Oracle RAC 10g Release 2/Oracle RAC 11g Release 1: The shared storage can be a cluster file system or raw VxVM volumes.
For Oracle RAC 11g Release 2: The shared storage can be a cluster file system or ASM disk groups created using raw VxVM volumes.
- Three or an odd number of standard disks or LUNs (recommended number is three) used as coordinator disks or as coordination point (CP) servers for I/O fencing.
- VCS manages the resources that are required by Oracle RAC. The resources must run in parallel on each node.

SF Oracle RAC includes the following technologies that are engineered to improve performance, availability, and manageability of Oracle RAC environments:

- Cluster File System (CFS) and Cluster Volume Manager (CVM) technologies to manage multi-instance database access to shared storage.
- An Oracle Disk Manager (ODM) library to maximize Oracle disk I/O performance.
- Interfaces to Oracle Clusterware/Grid Infrastructure and RAC for managing cluster membership and communication.

[Figure 1-2](#) displays the technologies that make up the SF Oracle RAC internal architecture.

Figure 1-2 SF Oracle RAC architecture



SF Oracle RAC provides an environment that can tolerate failures with minimal downtime and interruption to users. If a node fails as clients access the same database on multiple nodes, clients attached to the failed node can reconnect to a surviving node and resume access. Recovery after failure in the SF Oracle RAC environment is far quicker than recovery for a single-instance database because another Oracle instance is already up and running. The recovery process involves applying outstanding redo log entries of the failed node from the surviving nodes.

Component products and processes of SF Oracle RAC

[Table 1-1](#) lists the component products of SF Oracle RAC.

Table 1-1 SF Oracle RAC component products

Component product	Description
Cluster Volume Manager (CVM)	<p>Enables simultaneous access to shared volumes based on technology from Veritas Volume Manager (VxVM).</p> <p>See “Cluster Volume Manager (CVM)” on page 30.</p>
Cluster File System (CFS)	<p>Enables simultaneous access to shared file systems based on technology from Veritas File System (VxFS).</p> <p>See “Cluster File System (CFS)” on page 32.</p>
Veritas Cluster Server (VCS)	<p>Manages Oracle RAC databases and infrastructure components.</p> <p>See “Veritas Cluster Server” on page 35.</p>
Veritas I/O fencing	<p>Protects the data on shared disks when nodes in a cluster detect a change in the cluster membership that indicates a split-brain condition.</p> <p>See “About I/O fencing” on page 39.</p>
Oracle RAC	<p>Component of the Oracle database product that allows a database to be installed on multiple servers.</p> <p>See “Oracle RAC components” on page 81.</p>
Oracle Disk Manager (Database Accelerator)	<p>Provides the interface with the Oracle Disk Manager (ODM) API.</p> <p>See “Oracle Disk Manager” on page 84.</p>
RAC Extensions	<p>Manages cluster membership and communications between cluster nodes.</p> <p>See “RAC extensions” on page 85.</p>

Communication infrastructure

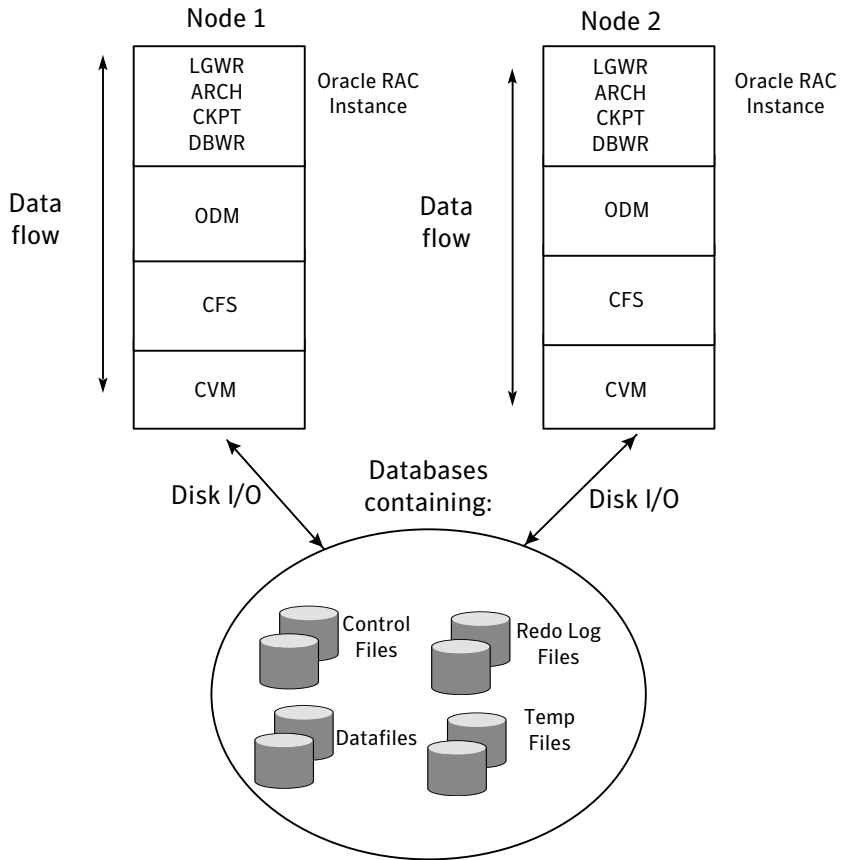
To understand the communication infrastructure, review the data flow and communication requirements.

Data flow

The CVM, CFS, ODM, and Oracle RAC elements reflect the overall data flow, or data stack, from an instance running on a server to the shared storage. The various Oracle processes composing an instance -- such as DBWR (Database Writer), LGWR (Log Writer process), CKPT (Storage Checkpoint process), and ARCH (Archive process (optional)) -- read and write data to the storage through the I/O stack. Oracle communicates through the ODM interface to CFS, which in turn accesses the storage through the CVM.

[Figure 1-3](#) represents the overall data flow.

Figure 1-3 Data stack

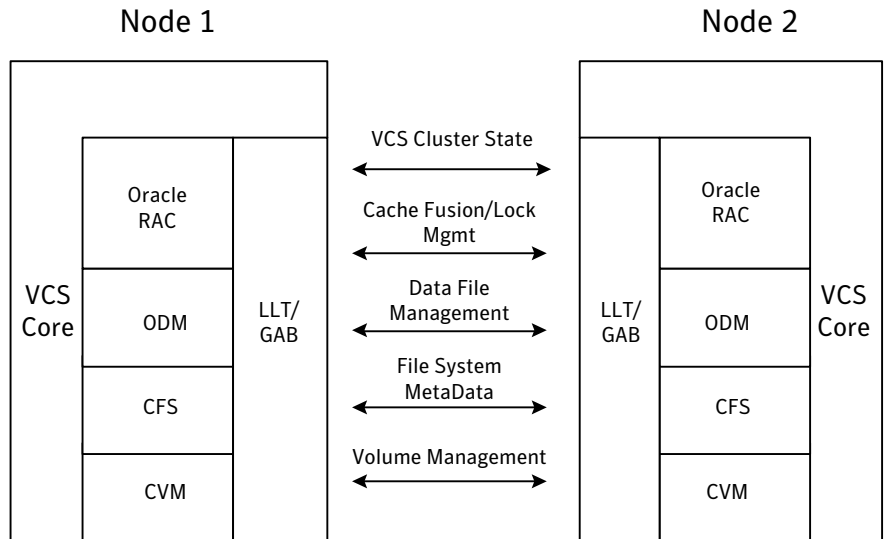


Communication requirements

End-users on a client system are unaware that they are accessing a database hosted by multiple instances. The key to performing I/O to a database accessed by multiple instances is communication between the processes. Each layer or component in the data stack must reliably communicate with its peer on other nodes to function properly. RAC instances must communicate to coordinate protection of data blocks in the database. ODM processes must communicate to coordinate data file protection and access across the cluster. CFS coordinates metadata updates for file systems, while CVM coordinates the status of logical volumes and maps.

Figure 1-4 represents the communication stack.

Figure 1-4 Communication stack



Cluster interconnect communication channel

The cluster interconnect provides an additional communication channel for all system-to-system communication, separate from the one-node communication between modules. Low Latency Transport (LLT) and Group Membership Services/Atomic Broadcast (GAB) make up the VCS communications package central to the operation of SF Oracle RAC.

In a standard operational state, significant traffic through LLT and GAB results from Lock Management, while traffic for other data is relatively sparse.

About Low Latency Transport (LLT)

The Low Latency Transport protocol is used for all cluster communications as a high-performance, low-latency replacement for the IP stack.

LLT has the following two major functions:

- Traffic distribution

LLT provides the communications backbone for GAB. LLT distributes (load balances) inter-system communication across all configured network links. This distribution ensures all cluster communications are evenly distributed across all network links for performance and fault resilience. If a link fails,

traffic is redirected to the remaining links. A maximum of eight network links are supported.

■ **Heartbeat**

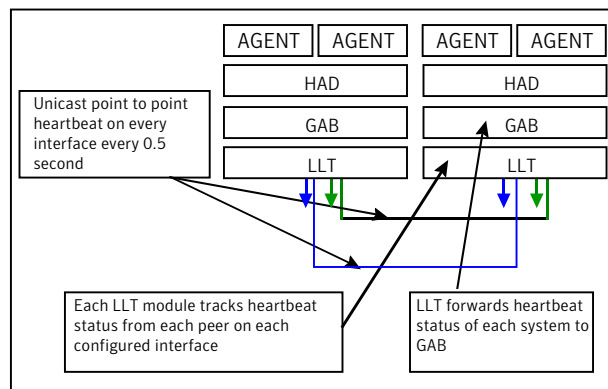
LLT is responsible for sending and receiving heartbeat traffic over each configured network link. The heartbeat traffic is point to point unicast. LLT uses ethernet broadcast to learn the address of the nodes in the cluster. All other cluster communications, including all status and configuration traffic is point to point unicast. The heartbeat is used by the Group Membership Services to determine cluster membership.

The heartbeat signal is defined as follows:

- LLT on each system in the cluster sends heartbeat packets out on all configured LLT interfaces every half second.
- LLT on each system tracks the heartbeat status from each peer on each configured LLT interface.
- LLT on each system forwards the heartbeat status of each system in the cluster to the local Group Membership Services function of GAB.
- GAB receives the status of heartbeat from all cluster systems from LLT and makes membership determination based on this information.

Figure 1-5 shows heartbeat in the cluster.

Figure 1-5 Heartbeat in the cluster



LLT can be configured to designate specific cluster interconnect links as either high priority or low priority. High priority links are used for cluster communications to GAB as well as heartbeat signals. Low priority links, during normal operation, are used for heartbeat and link state maintenance only, and the frequency of heartbeats is reduced to 50% of normal to reduce network overhead.

If there is a failure of all configured high priority links, LLT will switch all cluster communications traffic to the first available low priority link. Communication traffic will revert back to the high priority links as soon as they become available.

While not required, best practice recommends to configure at least one low priority link, and to configure two high priority links on dedicated cluster interconnects to provide redundancy in the communications path. Low priority links are typically configured on the public or administrative network.

If you use different media speed for the private NICs, Symantec recommends that you configure the NICs with lesser speed as low-priority links to enhance LLT performance. With this setting, LLT does active-passive load balancing across the private links. At the time of configuration and failover, LLT automatically chooses the link with high-priority as the active link and uses the low-priority links only when a high-priority link fails.

LLT sends packets on all the configured links in weighted round-robin manner. LLT uses the linkburst parameter which represents the number of back-to-back packets that LLT sends on a link before the next link is chosen. In addition to the default weighted round-robin based load balancing, LLT also provides destination-based load balancing. LLT implements destination-based load balancing where the LLT link is chosen based on the destination node id and the port. With destination-based load balancing, LLT sends all the packets of a particular destination on a link. However, a potential problem with the destination-based load balancing approach is that LLT may not fully utilize the available links if the ports have dissimilar traffic. Symantec recommends destination-based load balancing when the setup has more than two cluster nodes and more active LLT ports. You must manually configure destination-based load balancing for your cluster to set up the port to LLT link mapping.

See [“Configuring destination-based load balancing for LLT”](#) on page 122.

LLT on startup sends broadcast packets with LLT node id and cluster id information onto the LAN to discover any node in the network that has same node id and cluster id pair. Each node in the network replies to this broadcast message with its cluster id, node id, and node name.

LLT on the original node does not start and gives appropriate error in the following cases:

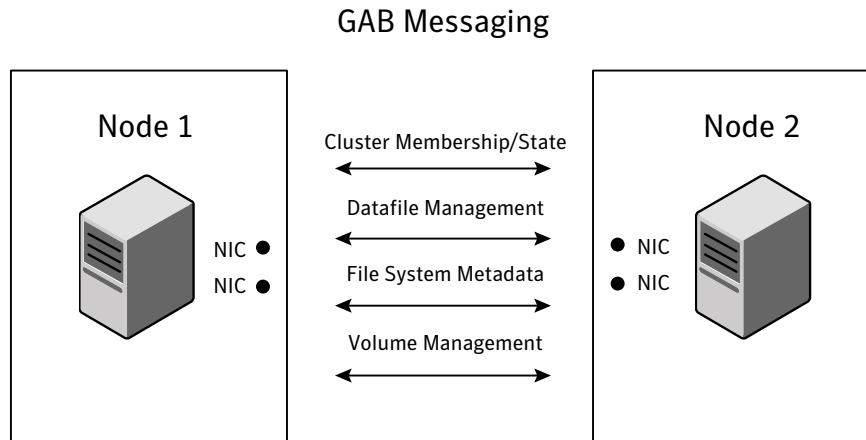
- LLT on any other node in the same network is running with the same node id and cluster id pair that it owns.
- LLT on the original node receives response from a node that does not have a node name entry in the `/etc/llthosts` file.

Group Membership Services/Atomic Broadcast

The GAB protocol is responsible for cluster membership and cluster communications.

Figure 1-6 shows the cluster communication using GAB messaging.

Figure 1-6 Cluster communication



Review the following information on cluster membership and cluster communication:

- Cluster membership

At a high level, all nodes configured by the installer can operate as a cluster; these nodes form a cluster membership. In SF Oracle RAC, a cluster membership specifically refers to all systems configured with the same cluster ID communicating by way of a redundant cluster interconnect.

All nodes in a distributed system, such as SF Oracle RAC, must remain constantly alert to the nodes currently participating in the cluster. Nodes can leave or join the cluster at any time because of shutting down, starting up, rebooting, powering off, or faulting processes. SF Oracle RAC uses its cluster membership capability to dynamically track the overall cluster topology.

SF Oracle RAC uses LLT heartbeats to determine cluster membership:

- When systems no longer receive heartbeat messages from a peer for a predetermined interval, a protocol excludes the peer from the current membership.
- GAB informs processes on the remaining nodes that the cluster membership has changed; this action initiates recovery actions specific to each module.

For example, CVM must initiate volume recovery and CFS must perform a fast parallel file system check.

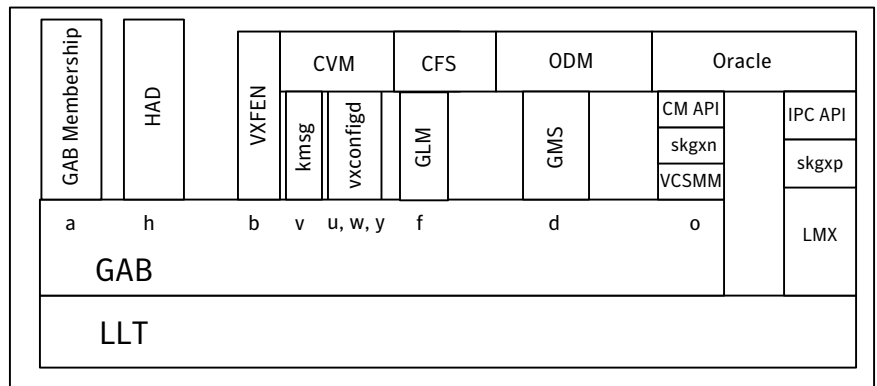
- When systems start receiving heartbeats from a peer outside of the current membership, a protocol enables the peer to join the membership.
- Cluster communications
 GAB provides reliable cluster communication between SF Oracle RAC modules. GAB provides guaranteed delivery of point-to-point messages and broadcast messages to all nodes. Point-to-point messaging involves sending and acknowledging the message. Atomic-broadcast messaging ensures all systems within the cluster receive all messages. If a failure occurs while transmitting a broadcast message, GAB ensures all systems have the same information after recovery.

Low-level communication: port relationship between GAB and processes

All components in SF Oracle RAC use GAB for communication. Each process that wants to communicate with a peer process on other nodes registers with GAB on a specific port. This registration enables communication and notification of membership changes. For example, the VCS engine (HAD) registers on port h. HAD receives messages from peer HAD processes on port h. HAD also receives notification when a node fails or when a peer process on port h unregisters.

Some modules use multiple ports for specific communications requirements. For example, CVM uses multiple ports to allow communications by kernel and user-level functions in CVM independently.

Figure 1-7 Low-level communication



Cluster Volume Manager (CVM)

CVM is an extension of Veritas Volume Manager, the industry-standard storage virtualization platform. CVM extends the concepts of VxVM across multiple nodes. Each node recognizes the same logical volume layout, and more importantly, the same state of all volume resources.

CVM supports performance-enhancing capabilities, such as striping, mirroring, and mirror break-off (snapshot) for off-host backup. You can use standard VxVM commands from one node in the cluster to manage all storage. All other nodes immediately recognize any changes in disk group and volume configuration with no user interaction.

For detailed information, see the *Veritas Storage Foundation Cluster File System High Availability Administrator's Guide*.

CVM architecture

CVM is designed with a "master and slave" architecture. One node in the cluster acts as the configuration master for logical volume management, and all other nodes are slaves. Any node can take over as master if the existing master fails. The CVM master exists on a per-cluster basis and uses GAB and LLT to transport its configuration data.

Just as with VxVM, the Volume Manager configuration daemon, `vxconfigd`, maintains the configuration of logical volumes. This daemon handles changes to the volumes by updating the operating system at the kernel level. For example, if a mirror of a volume fails, the mirror detaches from the volume and `vxconfigd` determines the proper course of action, updates the new volume layout, and informs the kernel of a new volume layout. CVM extends this behavior across multiple nodes and propagates volume changes to the master `vxconfigd`.

Note: You must perform operator-initiated changes on the master node.

The `vxconfigd` process on the master pushes these changes out to slave `vxconfigd` processes, each of which updates the local kernel. The kernel module for CVM is `kmsg`.

See [Figure 1-7](#) on page 29.

CVM does not impose any write locking between nodes. Each node is free to update any area of the storage. All data integrity is the responsibility of the upper application. From an application perspective, standalone systems access logical volumes in the same way as CVM systems.

By default, CVM imposes a "Uniform Shared Storage" model. All nodes must connect to the same disk sets for a given disk group. Any node unable to detect the entire set of physical disks for a given disk group cannot import the group. If a node loses contact with a specific disk, CVM excludes the node from participating in the use of that disk.

Set the `storage_connectivity` tunable to asymmetric to enable a cluster node to join even if the node does not have access to all of the shared storage. Similarly, a node can import a shared disk group even if there is a local failure to the storage.

For detailed information, see the *Veritas Storage Foundation Cluster File System High Availability Administrator's Guide*.

CVM communication

CVM communication involves various GAB ports for different types of communication. For an illustration of these ports:

See [Figure 1-7](#) on page 29.

CVM communication involves the following GAB ports:

- **Port w**

Most CVM communication uses port w for vxconfigd communications. During any change in volume configuration, such as volume creation, plex attachment or detachment, and volume resizing, vxconfigd on the master node uses port w to share this information with slave nodes.

When all slaves use port w to acknowledge the new configuration as the next active configuration, the master updates this record to the disk headers in the VxVM private region for the disk group as the next configuration.
- **Port v**

CVM uses port v for kernel-to-kernel communication. During specific configuration events, certain actions require coordination across all nodes. An example of synchronizing events is a resize operation. CVM must ensure all nodes see the new or old size, but never a mix of size among members.

CVM also uses this port to obtain cluster membership from GAB and determine the status of other CVM members in the cluster.
- **Port u**

CVM uses the group atomic broadcast (GAB) port u to ship the commands from the slave node to the master node.
- **Port y**

CVM uses port y for kernel-to-kernel communication required while shipping I/Os from nodes that might have lost local access to storage to other nodes in the cluster.

CVM recovery

When a node leaves a cluster, the new membership is delivered by GAB, to CVM on existing cluster nodes. The fencing driver (VXFEN) ensures that split-brain scenarios are taken care of before CVM is notified. CVM then initiates recovery of mirrors of shared volumes that might have been in an inconsistent state following the exit of the node.

For database files, when ODM is enabled with SmartSync option, Oracle Resilvering handles recovery of mirrored volumes. For non-database files, this recovery is optimized using Dirty Region Logging (DRL). The DRL is a map stored in a special purpose VxVM sub-disk and attached as an additional plex to the mirrored volume. When a DRL subdisk is created for a shared volume, the length of the sub-disk is automatically evaluated so as to cater to the number of cluster nodes. If the shared volume has Fast Mirror Resync (FlashSnap) enabled, the DCO (Data Change Object) log volume created automatically has DRL embedded in it. In the absence of DRL or DCO, CVM does a full mirror resynchronization.

Configuration differences with VxVM

CVM configuration differs from VxVM configuration in the following areas:

- Configuration commands occur on the master node.
- Disk groups are created and imported as shared disk groups. (Disk groups can also be private.)
- Disk groups are activated per node.
- Shared disk groups are automatically imported when CVM starts.

Cluster File System (CFS)

CFS enables you to simultaneously mount the same file system on multiple nodes and is an extension of the industry-standard Veritas File System. Unlike other file systems which send data through another node to the storage, CFS is a true SAN file system. All data traffic takes place over the storage area network (SAN), and only the metadata traverses the cluster interconnect.

In addition to using the SAN fabric for reading and writing data, CFS offers Storage Checkpoint and rollback for backup and recovery.

Access to cluster storage in typical SF Oracle RAC configurations use CFS. Raw access to CVM volumes is also possible but not part of a common configuration.

For detailed information, see the Veritas Storage Foundation Cluster File System High Availability Administrator's documentation.

CFS architecture

SF Oracle RAC uses CFS to manage a file system in a large database environment. Since CFS is an extension of VxFS, it operates in a similar fashion and caches metadata and data in memory (typically called buffer cache or vnode cache). CFS uses a distributed locking mechanism called Global Lock Manager (GLM) to ensure all nodes have a consistent view of the file system. GLM provides metadata and cache coherency across multiple nodes by coordinating access to file system metadata, such as inodes and free lists. The role of GLM is set on a per-file system basis to enable load balancing.

CFS involves a primary/secondary architecture. One of the nodes in the cluster is the primary node for a file system. Though any node can initiate an operation to create, delete, or resize data, the GLM master node carries out the actual operation. After creating a file, the GLM master node grants locks for data coherency across nodes. For example, if a node tries to modify a block in a file, it must obtain an exclusive lock to ensure other nodes that may have the same file cached have this cached copy invalidated.

SF Oracle RAC configurations minimize the use of GLM locking. Oracle RAC accesses the file system through the ODM interface and handles its own locking; only Oracle (and not GLM) buffers data and coordinates write operations to files. A single point of locking and buffering ensures maximum performance. GLM locking is only involved when metadata for a file changes, such as during create and resize operations.

CFS communication

CFS uses port `f` for GLM lock and metadata communication. SF Oracle RAC configurations minimize the use of GLM locking except when metadata for a file changes.

CFS file system benefits

Many features available in VxFS do not come into play in an SF Oracle RAC environment because ODM handles such features. CFS adds such features as high availability, consistency and scalability, and centralized management to VxFS. Using CFS in an SF Oracle RAC environment provides the following benefits:

- Increased manageability, including easy creation and expansion of files
In the absence of CFS, you must provide Oracle with fixed-size partitions. With CFS, you can grow file systems dynamically to meet future requirements.
- Less prone to user error

Raw partitions are not visible and administrators can compromise them by mistakenly putting file systems over the partitions. Nothing exists in Oracle to prevent you from making such a mistake.

- **Data center consistency**

If you have raw partitions, you are limited to a RAC-specific backup strategy. CFS enables you to implement your backup strategy across the data center.

CFS configuration differences

The first node to mount a CFS file system as shared becomes the primary node for that file system. All other nodes are "secondaries" for that file system.

Mount the cluster file system individually from each node. The `-o cluster` option of the `mount` command mounts the file system in shared mode, which means you can mount the file system simultaneously on mount points on multiple nodes.

When using the `fsadm` utility for online administration functions on VxFS file systems, including file system resizing, defragmentation, directory reorganization, and querying or changing the `largefiles` flag, run `fsadm` from any node.

CFS recovery

The `vxfsckd` daemon is responsible for ensuring file system consistency when a node crashes that was a primary node for a shared file system. If the local node is a secondary node for a given file system and a reconfiguration occurs in which this node becomes the primary node, the kernel requests `vxfsckd` on the new primary node to initiate a replay of the intent log of the underlying volume. The `vxfsckd` daemon forks a special call to `fsck` that ignores the volume reservation protection normally respected by `fsck` and other VxFS utilities. The `vxfsckd` can check several volumes at once if the node takes on the primary role for multiple file systems.

After a secondary node crash, no action is required to recover file system integrity. As with any crash on a file system, internal consistency of application data for applications running at the time of the crash is the responsibility of the applications.

Comparing raw volumes and CFS for data files

Keep these points in mind about raw volumes and CFS for data files:

- If you use file-system-based data files, the file systems containing these files must be located on shared disks. Create the same file system mount point on each node.

- If you use raw devices, such as VxVM volumes, set the permissions for the volumes to be owned permanently by the database account. VxVM sets volume permissions on import. The VxVM volume, and any file system that is created in it, must be owned by the Oracle database user.

Veritas Cluster Server

Veritas Cluster Server (VCS) directs SF Oracle RAC operations by controlling the startup and shutdown of components layers and providing monitoring and notification for failures.

In a typical SF Oracle RAC configuration, the Oracle RAC service groups for VCS run as "parallel" service groups rather than "failover" service groups; in the event of a failure, VCS does not attempt to migrate a failed service group. Instead, the software enables you to configure the group to restart on failure.

VCS architecture

The High Availability Daemon (HAD) is the main VCS daemon running on each node. HAD tracks changes in the cluster configuration and monitors resource status by communicating over GAB and LLT. HAD manages all application services using agents, which are installed programs to manage resources (specific hardware or software entities).

The VCS architecture is modular for extensibility and efficiency. HAD does not need to know how to start up Oracle or any other application under VCS control. Instead, you can add agents to manage different resources with no effect on the engine (HAD). Agents only communicate with HAD on the local node and HAD communicates status with HAD processes on other nodes. Because agents do not need to communicate across systems, VCS is able to minimize traffic on the cluster interconnect.

SF Oracle RAC provides specific agents for VCS to manage CVM, CFS, and Oracle components like Oracle Grid Infrastructure and database (including instances).

VCS communication

VCS uses port `h` for HAD communication. Agents communicate with HAD on the local node about resources, and HAD distributes its view of resources on that node to other nodes through GAB port `h`. HAD also receives information from other cluster members to update its own view of the cluster.

About the IMF notification module

The notification module of Intelligent Monitoring Framework (IMF) is the Asynchronous Monitoring Framework (AMF).

AMF is a kernel driver which hooks into system calls and other kernel interfaces of the operating system to get notifications on various events such as:

- When a process starts or stops.
- When a block device gets mounted or unmounted from a mount point.

AMF also interacts with the Intelligent Monitoring Framework Daemon (IMFD) to get disk group related notifications. AMF relays these notification to various VCS Agent that are enabled for intelligent monitoring.

See [“About resource monitoring”](#) on page 36.

About resource monitoring

VCS agents poll the resources periodically based on the monitor interval (in seconds) value that is defined in the MonitorInterval or in the OfflineMonitorInterval resource type attributes. After each monitor interval, VCS invokes the monitor agent function for that resource. For example, for process offline monitoring, the process agent's monitor agent function corresponding to each process resource scans the process table in each monitor interval to check whether the process has come online. For process online monitoring, the monitor agent function queries the operating system for the status of the process id that it is monitoring. In case of the mount agent, the monitor agent function corresponding to each mount resource checks if the block device is mounted on the mount point or not. In order to determine this, the monitor function does operations such as mount table scans or runs `statfs` equivalents.

With intelligent monitoring framework (IMF), VCS supports intelligent resource monitoring in addition to poll-based monitoring. IMF is an extension to the VCS agent framework. You can enable or disable the intelligent monitoring functionality of the VCS agents that are IMF-aware. For a list of IMF-aware agents, see the *Veritas Cluster Server Bundled Agents Reference Guide*.

See [“How intelligent resource monitoring works”](#) on page 37.

See [“Enabling and disabling intelligent resource monitoring for agents manually”](#) on page 122.

Poll-based monitoring can consume a fairly large percentage of system resources such as CPU and memory on systems with a huge number of resources. This not only affects the performance of running applications, but also places a limit on how many resources an agent can monitor efficiently.

However, with IMF-based monitoring you can either eliminate poll-based monitoring completely or reduce its frequency. For example, for process offline and online monitoring, you can completely avoid the need for poll-based monitoring with IMF-based monitoring enabled for processes. Similarly for vxfs

mounts, you can eliminate the poll-based monitoring with IMF monitoring enabled. Such reduction in monitor footprint will make more system resources available for other applications to consume.

Note: Intelligent Monitoring Framework for mounts is supported only for the VxFS, CFS, and NFS mount types.

With IMF-enabled agents, VCS will be able to effectively monitor larger number of resources.

Thus, intelligent monitoring has the following benefits over poll-based monitoring:

- Provides faster notification of resource state changes
- Reduces VCS system utilization due to reduced monitor function footprint
- Enables VCS to effectively monitor a large number of resources

Consider enabling IMF for an agent in the following cases:

- You have a large number of process resources or mount resources under VCS control.
- You have any of the agents that are IMF-aware.

How intelligent resource monitoring works

When an IMF-aware agent starts up, the agent initializes with the IMF notification module. After the resource moves to a steady state, the agent registers the details that are required to monitor the resource with the IMF notification module. For example, the process agent registers the PIDs of the processes with the IMF notification module. The agent's `imf_getnotification` function waits for any resource state changes. When the IMF notification module notifies the `imf_getnotification` function about a resource state change, the agent framework runs the monitor agent function to ascertain the state of that resource. The agent notifies the state change to VCS which takes appropriate action.

A resource moves into a steady state when any two consecutive monitor agent functions report the state as ONLINE or as OFFLINE. The following are a few examples of how steady state is reached.

- When a resource is brought online, a monitor agent function is scheduled after the online agent function is complete. Assume that this monitor agent function reports the state as ONLINE. The next monitor agent function runs after a time interval specified by the `MonitorInterval` attribute. The default value of `MonitorInterval` is 60 seconds. If this monitor agent function too reports the state as ONLINE, a steady state is achieved because two consecutive monitor agent functions reported the resource state as ONLINE. After the second

monitor agent function reports the state as ONLINE, the registration command for IMF is scheduled. The resource is registered with the IMF notification module and the resource comes under IMF control.

A similar sequence of events applies for taking a resource offline.

- When a resource is brought online, a monitor agent function is scheduled after the online agent function is complete. Assume that this monitor agent function reports the state as ONLINE. If you initiate a probe operation on the resource before the time interval specified by MonitorInterval, the probe operation invokes the monitor agent function immediately. If this monitor agent function again reports the state as ONLINE, a steady state is achieved because two consecutive monitor agent functions reported the resource state as ONLINE. After the second monitor agent function reports the state as ONLINE, the registration command for IMF is scheduled. The resource is registered with the IMF notification module and the resource comes under IMF control. A similar sequence of events applies for taking a resource offline.
- Assume that IMF is disabled for an agent type and you enable IMF for the agent type when the resource is ONLINE. The next monitor agent function occurs after a time interval specified by MonitorInterval. If this monitor agent function again reports the state as ONLINE, a steady state is achieved because two consecutive monitor agent functions reported the resource state as ONLINE. A similar sequence of events applies if the resource is OFFLINE initially and the next monitor agent function also reports the state as OFFLINE after you enable IMF for the agent type.
- Assume that IMF is disabled for an agent type and you enable IMF for the agent type when the resource is ONLINE. If you initiate a probe operation on the resource, this probe operation invokes the monitor agent function immediately. If this monitor agent function also reports the state as ONLINE, a steady state is achieved because two consecutive monitor agent functions reported the resource state as ONLINE. A similar sequence of events applies if the resource is OFFLINE initially and the next monitor agent function initiated by the probe operation also reports the state as OFFLINE after you enable IMF for the agent type.

See [“About the IMF notification module”](#) on page 35.

Cluster configuration files

VCS uses two configuration files in a default configuration:

- The main.cf file defines the entire cluster, including the cluster name, systems in the cluster, and definitions of service groups and resources, in addition to service group and resource dependencies.

- The `types.cf` file defines the resource types. Each resource in a cluster is identified by a unique name and classified according to its type. VCS includes a set of pre-defined resource types for storage, networking, and application services.

Additional files similar to `types.cf` may be present if you add agents. For example, SF Oracle RAC includes additional resource types files, such as `OracleTypes.cf`, `PrivNIC.cf`, and `MultiPrivNIC.cf`.

About I/O fencing

I/O fencing protects the data on shared disks when nodes in a cluster detect a change in the cluster membership that indicates a split-brain condition.

The fencing operation determines the following:

- The nodes that must retain access to the shared storage
- The nodes that must be ejected from the cluster

This decision prevents possible data corruption. The installer installs the I/O fencing driver, `VRTSvxfen` depot, when you install SF Oracle RAC. To protect data on shared disks, you must configure I/O fencing after you install and configure SF Oracle RAC.

I/O fencing technology uses coordination points for arbitration in the event of a network partition.

I/O fencing coordination points can be coordinator disks or coordination point servers (CP servers) or both. You can configure disk-based or server-based I/O fencing:

Disk-based I/O fencing

I/O fencing that uses coordinator disks is referred to as disk-based I/O fencing.

Disk-based I/O fencing ensures data integrity in a single cluster.

Server-based I/O fencing

I/O fencing that uses at least one CP server system is referred to as server-based I/O fencing.

Server-based fencing can include only CP servers, or a mix of CP servers and coordinator disks.

Server-based I/O fencing ensures data integrity in clusters.

For detailed information, see the *Veritas Cluster Server Administrator's Guide*.

See [“About preventing data corruption with I/O fencing”](#) on page 44.

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide*.

About server-based I/O fencing

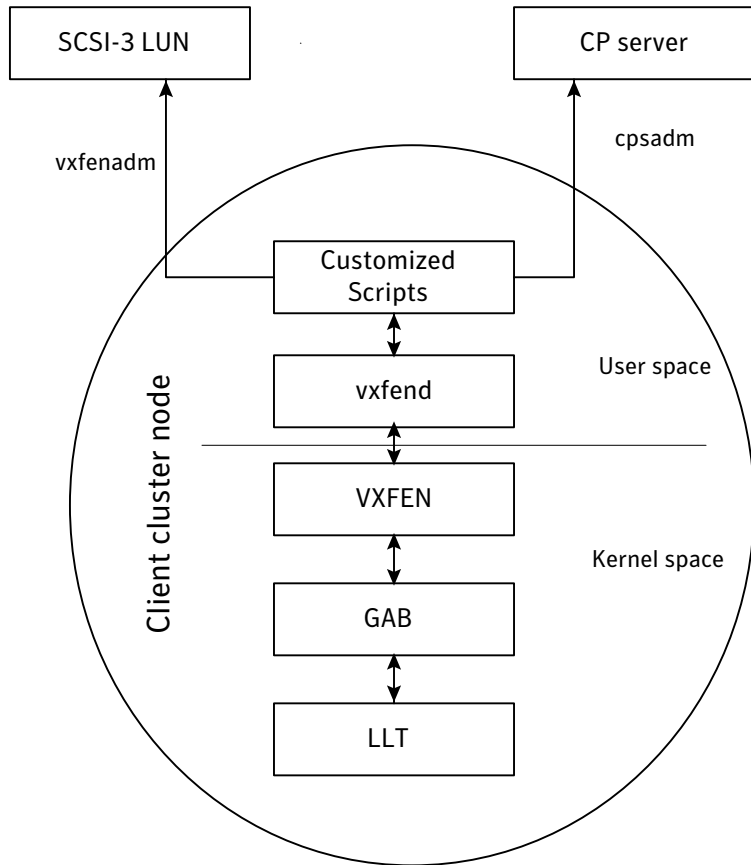
In a disk-based I/O fencing implementation, the vxfen driver handles various SCSI-3 PR based arbitration operations completely within the driver. I/O fencing also provides a framework referred to as customized fencing wherein arbitration operations are implemented in custom scripts. The vxfen driver invokes the custom scripts.

The CP server-based coordination point uses a customized fencing framework. Note that SCSI-3 PR based fencing arbitration can also be enabled using customized fencing framework. This allows the user to specify a combination of SCSI-3 LUNs and CP servers as coordination points using customized fencing. Customized fencing can be enabled by specifying `vxfen_mode=customized` and `vxfen_mechanism=cps` in the `/etc/vxfenmode` file.

Moreover, both `/etc/vxfenmode` and `/etc/vxfentab` files contain additional parameter "security" which indicates if communication between CP server and SF Oracle RAC cluster nodes is secure.

[Figure 1-8](#) displays a schematic of the customized fencing options.

Figure 1-8 Customized fencing



A user level daemon vxfend interacts with the vxfen driver, which in turn interacts with GAB to get the node membership update. Upon receiving membership updates, vxfend invokes various scripts to race for the coordination point and fence off data disks. The vxfend daemon manages various fencing agents. The customized fencing scripts are located in the `/opt/VRTSvcs/vxfen/bin/customized/cps` directory.

Table 1-2 The scripts that are involved include the following:

Script name	Description
generate_snapshot.sh	Retrieves the SCSI ID's of the coordinator disks and/or UUID ID's of the CP servers CP server uses the UUID stored in <code>/etc/VRTScps/db/current/cps_uuid</code> . For information about the UUID (Universally Unique Identifier), see the <i>Veritas Cluster Server Administrator's Guide</i> .
join_local_node.sh	Registers the keys with the coordinator disks or CP servers.
race_for_coordination_point.sh:	Races to determine a winner after cluster reconfiguration.
unjoin_local_node.sh	Removes the keys that are registered in <code>join_local_node.sh</code> .
fence_data_disks.sh	Fences the data disks from access by the losing nodes.
local_info.sh:	Lists local node's configuration parameters and coordination points, which are used by the <code>vxfen</code> driver.

I/O fencing enhancements provided by CP server

CP server configurations enhance disk-based I/O fencing by providing the following new capabilities:

- CP server configurations are scalable, and a configuration with three CP servers can provide I/O fencing for multiple SF Oracle RAC clusters. Since a single CP server configuration can serve a large number of SF Oracle RAC clusters, the cost of multiple SF Oracle RAC cluster deployments can be significantly reduced.
- Appropriately situated CP servers can eliminate any coordinator disk location bias in the I/O fencing process. For example, this location bias may occur where, due to logistical restrictions, two of the three coordinator disks are located at a single site, and the cost of setting up a third coordinator disk location is prohibitive.

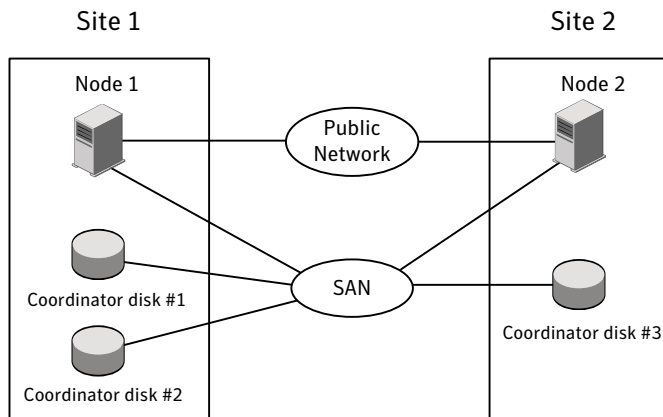
See [Figure 1-9](#) on page 43.

In such a configuration, if the site with two coordinator disks is inaccessible, the other site does not survive due to a lack of a majority of coordination points.

I/O fencing would require extension of the SAN to the third site which may not be a suitable solution. An alternative is to place a CP server at a remote site as the third coordination point.

Note: The CP server provides an alternative arbitration mechanism without having to depend on SCSI-3 compliant coordinator disks. Data disk fencing in Cluster Volume Manager (CVM) will still require SCSI-3 I/O fencing.

Figure 1-9 Skewed placement of coordinator disks at Site 1



About the CP server database

CP server requires a database for storing the registration keys of the SF Oracle RAC cluster nodes. CP server uses a SQLite database for its operations. By default, the database is located at `/etc/VRTScps/db`.

For a single node VCS cluster hosting a CP server, the database can be placed on a local file system. For an SFHA cluster hosting a CP server, the database must be placed on a shared file system. The file system must be shared among all nodes that are part of the SFHA cluster.

In an SFHA cluster hosting the CP server, the shared database is protected by setting up SCSI-3 PR based I/O fencing. SCSI-3 PR based I/O fencing protects against split-brain scenarios.

Warning: The CP server database must not be edited directly and should only be accessed using `cpsadm(1M)`. Manipulating the database manually may lead to undesirable results including system panics.

About the CP server user types and privileges

The CP server supports the following user types, each with a different access level privilege:

- CP server administrator (admin)
- CP server operator

Different access level privileges permit the user to issue different commands. If a user is neither a CP server admin nor a CP server operator user, then the user has guest status and can issue limited commands.

The user types and their access level privileges are assigned to individual users during SF Oracle RAC cluster configuration for fencing. During the installation process, you are prompted for a user name, password, and access level privilege (CP server admin or CP server operator).

To administer and operate a CP server, there must be at least one CP server admin.

A root user on a CP server is given all the administrator privileges, and these administrator privileges can be used to perform all the CP server specific operations.

About preferred fencing

The I/O fencing driver uses coordination points to prevent split-brain in a VCS cluster. By default, the fencing driver favors the subcluster with maximum number of nodes during the race for coordination points. With the preferred fencing feature, you can specify how the fencing driver must determine the surviving subcluster.

You can configure the preferred fencing policy using the cluster-level attribute PreferredFencingPolicy for the following:

- Enable system-based preferred fencing policy to give preference to high capacity systems.
- Enable group-based preferred fencing policy to give preference to service groups for high priority applications.
- Disable preferred fencing policy to use the default node count-based race policy.

See [“Enabling or disabling the preferred fencing policy”](#) on page 154.

About preventing data corruption with I/O fencing

I/O fencing is a feature that prevents data corruption in the event of a communication breakdown in a cluster.

To provide high availability, the cluster must be capable of taking corrective action when a node fails. In this situation, SF Oracle RAC configures its components to reflect the altered membership.

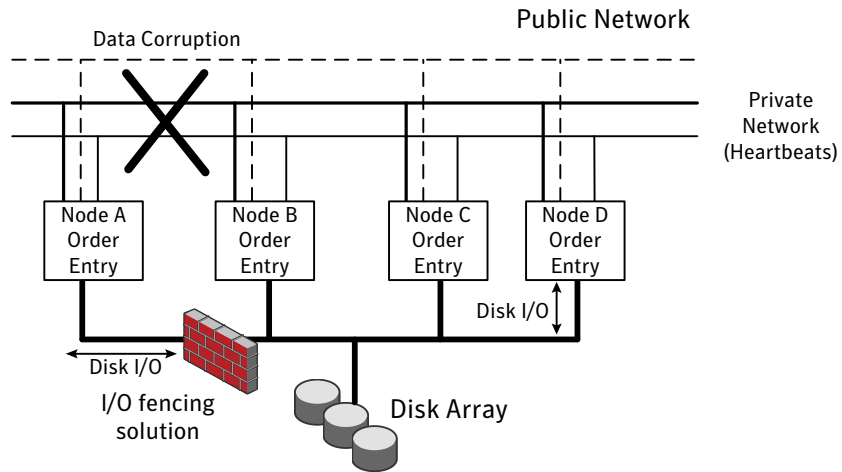
Problems arise when the mechanism that detects the failure breaks down because symptoms appear identical to those of a failed node. For example, if a system in a two-node cluster fails, the system stops sending heartbeats over the private interconnects. The remaining node then takes corrective action. The failure of the private interconnects, instead of the actual nodes, presents identical symptoms and causes each node to determine its peer has departed. This situation typically results in data corruption because both nodes try to take control of data storage in an uncoordinated manner.

In addition to a broken set of private networks, other scenarios can generate this situation. If a system is so busy that it appears to stop responding or "hang," the other nodes could declare it as dead. This declaration may also occur for the nodes that use the hardware that supports a "break" and "resume" function. When a node drops to PROM level with a break and subsequently resumes operations, the other nodes may declare the system dead. They can declare it dead even if the system later returns and begins write operations.

SF Oracle RAC uses I/O fencing to remove the risk that is associated with split-brain. I/O fencing allows write access for members of the active cluster. It blocks access to storage from non-members.

[Figure 1-10](#) displays a schematic of a four node cluster, each node writing order entries to the connected disk array. When the private network connection between the four nodes is disrupted (between Node A and the other 3 nodes in the figure below), a split-brain situation occurs with the possibility of data corruption to the disk array. The I/O fencing process prevents split-brain and any data corruption by fencing off Node A from the cluster.

Figure 1-10 Private network disruption and I/O fencing solution



About SCSI-3 Persistent Reservations

SCSI-3 Persistent Reservations (SCSI-3 PR) are required for I/O fencing and resolve the issues of using SCSI reservations in a clustered SAN environment. SCSI-3 PR enables access for multiple nodes to a device and simultaneously blocks access for other nodes.

SCSI-3 reservations are persistent across SCSI bus resets and support multiple paths from a host to a disk. In contrast, only one host can use SCSI-2 reservations with one path. If the need arises to block access to a device because of data integrity concerns, only one host and one path remain active. The requirements for larger clusters, with multiple nodes reading and writing to storage in a controlled manner, make SCSI-2 reservations obsolete.

SCSI-3 PR uses a concept of registration and reservation. Each system registers its own "key" with a SCSI-3 device. Multiple systems registering keys form a membership and establish a reservation, typically set to "Write Exclusive Registrants Only (WERO)." The WERO setting enables only registered systems to perform write operations. For a given disk, only one reservation can exist amidst numerous registrations.

With SCSI-3 PR technology, blocking write access is as easy as removing a registration from a device. Only registered members can "eject" the registration of another member. A member wishing to eject another member issues a "preempt and abort" command. Ejecting a node is final and atomic; an ejected node cannot eject another node. In SF Oracle RAC, a node registers the same key for all paths to the device. A single preempt and abort command ejects a node from all paths to the storage device.

About I/O fencing operations

I/O fencing, provided by the kernel-based fencing module (`vxxfen`), performs identically on node failures and communications failures. When the fencing module on a node is informed of a change in cluster membership by the GAB module, it immediately begins the fencing operation. The node tries to eject the key for departed nodes from the coordinator disks using the `preempt` and `abort` command. When the node successfully ejects the departed nodes from the coordinator disks, it also ejects the departed nodes from the data disks. In a split-brain scenario, both sides of the split would race for control of the coordinator disks. The side winning the majority of the coordinator disks wins the race and fences the loser. The loser then panics and restarts the system.

See [“About I/O fencing components”](#) on page 47.

See [“How I/O fencing works in different event scenarios”](#) on page 49.

About I/O fencing communication

The `vxxfen` driver connects to GAB port `b` to intercept cluster membership changes (reconfiguration messages). During a membership change, the fencing driver determines which systems are members of the cluster to allow access to shared disks.

After completing fencing operations, the driver passes reconfiguration messages to higher modules. CVM handles fencing of data drives for shared disk groups. After a node successfully joins the GAB cluster and the driver determines that a preexisting split-brain does not exist, CVM can import all shared disk groups. The CVM master coordinates the order of import and the key for each disk group. As each slave joins the cluster, it accepts the CVM list of disk groups and keys, and adds its proper digit to the first byte of the key. Each slave then registers the keys with all drives in the disk groups.

About I/O fencing components

The shared storage for SF Oracle RAC must support SCSI-3 persistent reservations to enable I/O fencing. SF Oracle RAC involves two types of shared storage:

- Data disks—Store shared data
See [“About data disks”](#) on page 47.
- Coordination points—Act as a global lock during membership changes
See [“About coordination points”](#) on page 48.

About data disks

Data disks are standard disk devices for data storage and are either physical disks or RAID Logical Units (LUNs).

These disks must support SCSI-3 PR and must be part of standard VxVM or CVM disk groups. CVM is responsible for fencing data disks on a disk group basis. Disks that are added to a disk group and new paths that are discovered for a device are automatically fenced.

About coordination points

Coordination points provide a lock mechanism to determine which nodes get to fence off data drives from other nodes. A node must eject a peer from the coordination points before it can fence the peer from the data drives. SF Oracle RAC prevents split-brain when vxfen races for control of the coordination points and the winner partition fences the ejected nodes from accessing the data disks.

The coordination points can either be disks or servers or both.

■ Coordinator disks

Disks that act as coordination points are called coordinator disks. Coordinator disks are three standard disks or LUNs set aside for I/O fencing during cluster reconfiguration. Coordinator disks do not serve any other storage purpose in the SF Oracle RAC configuration.

Dynamic Multi-pathing (DMP) allows coordinator disks to take advantage of the path failover and the dynamic adding and removal capabilities of DMP.

On cluster nodes with HP-UX 11i v3, you must use DMP devices or iSCSI devices for I/O fencing. The following changes in HP-UX 11i v3 require you to not use raw devices for I/O fencing:

■ Provides native multi-pathing support

■ Does not provide access to individual paths through the device file entries

The metanode interface that HP-UX provides does not meet the SCSI-3 PR requirements for the I/O fencing feature. You can configure coordinator disks to use Veritas Volume Manager Dynamic Multi-pathing (DMP) feature.

See the *Veritas Storage Foundation Administrator's Guide*.

■ Coordination point servers

The coordination point server (CP server) is a software solution which runs on a remote system or cluster. CP server provides arbitration functionality by allowing the SF Oracle RAC cluster nodes to perform the following tasks:

■ Self-register to become a member of an active SF Oracle RAC cluster (registered with CP server) with access to the data drives

■ Check which other nodes are registered as members of this active SF Oracle RAC cluster

■ Self-unregister from this active SF Oracle RAC cluster

- Forcefully unregister other nodes (preempt) as members of this active SF Oracle RAC cluster

In short, the CP server functions as another arbitration mechanism that integrates within the existing I/O fencing module.

Note: With the CP server, the fencing arbitration logic still remains on the SF Oracle RAC cluster.

Multiple SF Oracle RAC clusters running different operating systems can simultaneously access the CP server. TCP/IP based communication is used between the CP server and the SF Oracle RAC clusters.

How I/O fencing works in different event scenarios

[Table 1-3](#) describes how I/O fencing works to prevent data corruption in different failure event scenarios. For each event, review the corrective operator actions.

Table 1-3 I/O fencing scenarios

Event	Node A: What happens?	Node B: What happens?	Operator action
Both private networks fail.	Node A races for majority of coordination points. If Node A wins race for coordination points, Node A ejects Node B from the shared disks and continues.	Node B races for majority of coordination points. If Node B loses the race for the coordination points, Node B panics and removes itself from the cluster.	When Node B is ejected from cluster, repair the private networks before attempting to bring Node B back.
Both private networks function again after event above.	Node A continues to work.	Node B has crashed. It cannot start the database since it is unable to write to the data disks.	Restart Node B after private networks are restored.
One private network fails.	Node A prints message about an IOFENCE on the console but continues.	Node B prints message about an IOFENCE on the console but continues.	Repair private network. After network is repaired, both nodes automatically use it.

Table 1-3 I/O fencing scenarios (*continued*)

Event	Node A: What happens?	Node B: What happens?	Operator action
Node A hangs.	<p>Node A is extremely busy for some reason or is in the kernel debugger.</p> <p>When Node A is no longer hung or in the kernel debugger, any queued writes to the data disks fail because Node A is ejected. When Node A receives message from GAB about being ejected, it panics and removes itself from the cluster.</p>	<p>Node B loses heartbeats with Node A, and races for a majority of coordination points.</p> <p>Node B wins race for coordination points and ejects Node A from shared data disks.</p>	Repair or debug the node that hangs and reboot the node to rejoin the cluster.

Table 1-3 I/O fencing scenarios (*continued*)

Event	Node A: What happens?	Node B: What happens?	Operator action
<p>Nodes A and B and private networks lose power. Coordination points and data disks retain power.</p> <p>Power returns to nodes and they restart, but private networks still have no power.</p>	<p>Node A restarts and I/O fencing driver (vxfen) detects Node B is registered with coordination points. The driver does not see Node B listed as member of cluster because private networks are down. This causes the I/O fencing device driver to prevent Node A from joining the cluster. Node A console displays:</p> <p>Potentially a preexisting split brain. Dropping out of the cluster. Refer to the user documentation for steps required to clear preexisting split brain.</p>	<p>Node B restarts and I/O fencing driver (vxfen) detects Node A is registered with coordination points. The driver does not see Node A listed as member of cluster because private networks are down. This causes the I/O fencing device driver to prevent Node B from joining the cluster. Node B console displays:</p> <p>Potentially a preexisting split brain. Dropping out of the cluster. Refer to the user documentation for steps required to clear preexisting split brain.</p>	<p>Resolve preexisting split-brain condition.</p> <p>See “Fencing startup reports preexisting split-brain” on page 202.</p>

Table 1-3 I/O fencing scenarios (*continued*)

Event	Node A: What happens?	Node B: What happens?	Operator action
<p>Node A crashes while Node B is down. Node B comes up and Node A is still down.</p>	<p>Node A is crashed.</p>	<p>Node B restarts and detects Node A is registered with the coordination points. The driver does not see Node A listed as member of the cluster. The I/O fencing device driver prints message on console:</p> <p>Potentially a preexisting split brain. Dropping out of the cluster. Refer to the user documentation for steps required to clear preexisting split brain.</p>	<p>Resolve preexisting split-brain condition.</p> <p>See “Fencing startup reports preexisting split-brain” on page 202.</p>
<p>The disk array containing two of the three coordination points is powered off.</p> <p>No node leaves the cluster membership</p>	<p>Node A continues to operate as long as no nodes leave the cluster.</p>	<p>Node B continues to operate as long as no nodes leave the cluster.</p>	<p>Power on the failed disk array so that subsequent network partition does not cause cluster shutdown, or replace coordination points.</p> <p>See “Replacing I/O fencing coordinator disks when the cluster is online” on page 143.</p>

Table 1-3 I/O fencing scenarios (*continued*)

Event	Node A: What happens?	Node B: What happens?	Operator action
<p>The disk array containing two of the three coordination points is powered off.</p> <p>Node B gracefully leaves the cluster and the disk array is still powered off. Leaving gracefully implies a clean shutdown so that vxfen is properly unconfigured.</p>	<p>Node A continues to operate in the cluster.</p>	<p>Node B has left the cluster.</p>	<p>Power on the failed disk array so that subsequent network partition does not cause cluster shutdown, or replace coordination points.</p> <p>See “Replacing I/O fencing coordinator disks when the cluster is online” on page 143.</p>
<p>The disk array containing two of the three coordination points is powered off.</p> <p>Node B abruptly crashes or a network partition occurs between node A and node B, and the disk array is still powered off.</p>	<p>Node A races for a majority of coordination points. Node A fails because only one of the three coordination points is available. Node A panics and removes itself from the cluster.</p>	<p>Node B has left cluster due to crash or network partition.</p>	<p>Power on the failed disk array and restart I/O fencing driver to enable Node A to register with all coordination points, or replace coordination points.</p> <p>See “Replacing defective disks when the cluster is offline” on page 206.</p>

About CP server

This section discusses the CP server features.

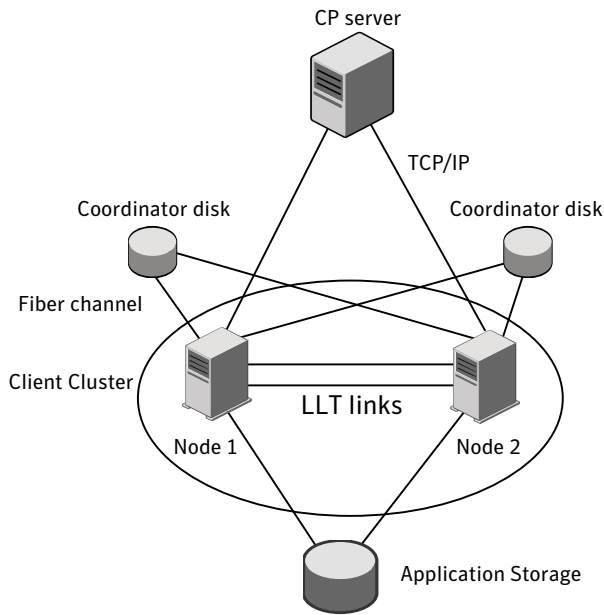
The following CP server features are described:

- SF Oracle RAC cluster configurations with server-based I/O fencing
- I/O fencing enhancements provided by the CP server
- About making CP server highly available
- Recommended CP server configurations
- About secure communication between the SF Oracle RAC cluster and CP server

Typical SF Oracle RAC cluster configuration with server-based I/O fencing

Figure 1-11 displays a configuration using a SF Oracle RAC cluster (with two nodes), a single CP server, and two coordinator disks. The nodes within the SF Oracle RAC cluster are connected to and communicate with each other using LLT links.

Figure 1-11 CP server, SF Oracle RAC cluster, and coordinator disks



Defining Coordination Points

Three or more odd number of coordination points are required for I/O fencing. A coordination point can be either a CP server or a coordinator disk. A CP server provides the same functionality as a coordinator disk in an I/O fencing scenario. Therefore, it is possible to mix and match CP servers and coordinator disks for the purpose of providing arbitration.

Symantec supports the following three coordination point configurations:

- Vxfen driver based I/O fencing using SCSI-3 coordinator disks
- Customized fencing using a combination of SCSI-3 disks and CP server(s) as coordination points
- Customized fencing using only three CP servers as coordination points

Note: Symantec does not support a configuration where multiple CP servers are configured on the same machine.

Deployment and migration scenarios for CP server

[Table 1-4](#) describes the supported deployment and migration scenarios, and the procedures you must perform on the SF Oracle RAC cluster and the CP server.

Table 1-4 CP server deployment and migration scenarios

Scenario	CP server	SF Oracle RAC cluster	Action required
Setup of CP server for a SF Oracle RAC cluster for the first time	New CP server	New SF Oracle RAC cluster using CP server as coordination point	<p>On the designated CP server, perform the following tasks:</p> <ol style="list-style-type: none"> 1 Prepare to configure the new CP server. 2 Configure the new CP server. 3 Prepare the new CP server for use by the SF Oracle RAC cluster. <p>On the SF Oracle RAC cluster nodes, configure server-based I/O fencing.</p> <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p>
Add a new SF Oracle RAC cluster to an existing and operational CP server	Existing and operational CP server	New SF Oracle RAC cluster	<p>On the SF Oracle RAC cluster nodes, configure server-based I/O fencing.</p> <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p>
Replace the coordination point from an existing CP server to a new CP server	New CP server	Existing SF Oracle RAC cluster using CP server as coordination point	<p>On the designated CP server, perform the following tasks:</p> <ol style="list-style-type: none"> 1 Prepare to configure the new CP server. 2 Configure the new CP server. 3 Prepare the new CP server for use by the SF Oracle RAC cluster. <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p> <p>On a node in the SF Oracle RAC cluster, run the <code>vxfsnwap</code> command to move to replace the CP server:</p> <p>See “Replacing coordination points for server-based fencing in an online cluster” on page 158.</p>

Table 1-4 CP server deployment and migration scenarios (*continued*)

Scenario	CP server	SF Oracle RAC cluster	Action required
Replace the coordination point from an existing CP server to an operational CP server coordination point	Operational CP server	Existing SF Oracle RAC cluster using CP server as coordination point	<p>On the designated CP server, prepare to configure the new CP server manually.</p> <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p> <p>On a node in the SF Oracle RAC cluster, run the <code>vxfsenwap</code> command to move to replace the CP server:</p> <p>See “Replacing coordination points for server-based fencing in an online cluster” on page 158.</p>
Enabling fencing in a SF Oracle RAC cluster with a new CP server coordination point	New CP server	Existing SF Oracle RAC cluster with fencing configured in disabled mode	<p>Note: Migrating from fencing in disabled mode to customized mode incurs application downtime on the SF Oracle RAC cluster.</p> <p>On the designated CP server, perform the following tasks:</p> <ol style="list-style-type: none"> 1 Prepare to configure the new CP server. 2 Configure the new CP server 3 Prepare the new CP server for use by the SF Oracle RAC cluster <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p> <p>On the SF Oracle RAC cluster nodes, perform the following:</p> <ol style="list-style-type: none"> 1 Stop all applications, VCS, and fencing on the SF Oracle RAC cluster. 2 To stop VCS, use the following command (to be run on all the SF Oracle RAC cluster nodes): <pre># hstop -local</pre> 3 Stop fencing using the following command: <pre># /sbin/init.d/vxfen stop</pre> 4 Reconfigure I/O fencing on the SF Oracle RAC cluster. <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p>

Table 1-4 CP server deployment and migration scenarios (*continued*)

Scenario	CP server	SF Oracle RAC cluster	Action required
<p>Enabling fencing in a SF Oracle RAC cluster with an operational CP server coordination point</p>	<p>Operational CP server</p>	<p>Existing SF Oracle RAC cluster with fencing configured in disabled mode</p>	<p>Note: Migrating from fencing in disabled mode to customized mode incurs application downtime.</p> <p>On the designated CP server, prepare to configure the new CP server.</p> <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for this procedure.</p> <p>On the SF Oracle RAC cluster nodes, perform the following tasks:</p> <ol style="list-style-type: none"> 1 Stop all applications, VCS, and fencing on the SF Oracle RAC cluster. 2 To stop VCS, use the following command (to be run on all the SF Oracle RAC cluster nodes): <pre># hstop -local</pre> 3 Stop fencing using the following command: <pre># /sbin/init.d/vxfen stop</pre> 4 Reconfigure fencing on the SF Oracle RAC cluster. <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p>

Table 1-4 CP server deployment and migration scenarios (*continued*)

Scenario	CP server	SF Oracle RAC cluster	Action required
<p>Enabling fencing in a SF Oracle RAC cluster with a new CP server coordination point</p>	<p>New CP server</p>	<p>Existing SF Oracle RAC cluster with fencing configured in scsi3 mode</p>	<p>On the designated CP server, perform the following tasks:</p> <ol style="list-style-type: none"> 1 Prepare to configure the new CP server. 2 Configure the new CP server 3 Prepare the new CP server for use by the SF Oracle RAC cluster <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p> <p>Based on whether the cluster is online or offline, perform the following procedures:</p> <p>For a cluster that is online, perform the following task on the SF Oracle RAC cluster:</p> <ul style="list-style-type: none"> ◆ Run the <code>vx fenceswap</code> command to migrate from disk-based fencing to the server-based fencing. <p>See “Migrating from disk-based to server-based fencing in an online cluster” on page 60.</p> <p>For a cluster that is offline, perform the following tasks on the SF Oracle RAC cluster:</p> <ol style="list-style-type: none"> 1 Stop all applications, VCS, and fencing on the SF Oracle RAC cluster. 2 To stop VCS, use the following command (to be run on all the SF Oracle RAC cluster nodes): <ul style="list-style-type: none"> # <code>hastop -local</code> 3 Stop fencing using the following command: <ul style="list-style-type: none"> # <code>/sbin/init.d/vxfen stop</code> 4 Reconfigure I/O fencing on the SF Oracle RAC cluster. <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p>

Table 1-4 CP server deployment and migration scenarios (*continued*)

Scenario	CP server	SF Oracle RAC cluster	Action required
<p>Enabling fencing in a SF Oracle RAC cluster with an operational CP server coordination point</p>	<p>Operational CP server</p>	<p>Existing SF Oracle RAC cluster with fencing configured in disabled mode</p>	<p>On the designated CP server, prepare to configure the new CP server.</p> <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for this procedure.</p> <p>Based on whether the cluster is online or offline, perform the following procedures:</p> <p>For a cluster that is online, perform the following task on the SF Oracle RAC cluster:</p> <ul style="list-style-type: none"> ◆ Run the <code>vx fenceswap</code> command to migrate from disk-based fencing to the server-based fencing. See “Migrating from disk-based to server-based fencing in an online cluster” on page 60. <p>For a cluster that is offline, perform the following tasks on the SF Oracle RAC cluster:</p> <ol style="list-style-type: none"> 1 Stop all applications, VCS, and fencing on the SF Oracle RAC cluster. 2 To stop VCS, use the following command (to be run on all the SF Oracle RAC cluster nodes): <code># hastop -local</code> 3 Stop fencing using the following command: <code># /sbin/init.d/vxfen stop</code> 4 Reconfigure fencing on the SF Oracle RAC cluster. <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> for the procedures.</p>

Table 1-4 CP server deployment and migration scenarios (*continued*)

Scenario	CP server	SF Oracle RAC cluster	Action required
Refreshing registrations of SF Oracle RAC cluster nodes on coordination points (CP servers/ coordinator disks) without incurring application downtime	Operational CP server	Existing SF Oracle RAC cluster using the CP server as coordination point	On the SF Oracle RAC cluster run the <code>vxfsenswap</code> command to refresh the keys on the CP server: See “Refreshing registration keys on the coordination points for server-based fencing” on page 156.

About migrating between disk-based and server-based fencing configurations

You can migrate between fencing configurations without incurring application downtime in the SF Oracle RAC clusters.

You can migrate from disk-based fencing to server-based fencing in the following cases:

- You want to leverage the benefits of server-based fencing.
- You want to replace faulty coordinator disks with coordination point servers (CP servers).

See [“Migrating from disk-based to server-based fencing in an online cluster”](#) on page 60.

Similarly, you can migrate from server-based fencing to disk-based fencing when you want to perform maintenance tasks on the CP server systems.

See [“Migrating from server-based to disk-based fencing in an online cluster”](#) on page 65.

Migrating from disk-based to server-based fencing in an online cluster

You can either use the installer or manually migrate from disk-based fencing to server-based fencing without incurring application downtime in the SF Oracle RAC clusters.

See [“About migrating between disk-based and server-based fencing configurations”](#) on page 60.

You can also use response files to migrate between fencing configurations.

See [“Migrating between fencing configurations using response files”](#) on page 71.

Warning: The cluster might panic if any node leaves the cluster membership before the coordination points migration operation completes.

This section covers the following procedures:

- | | |
|--|--|
| Migrating using the script-based installer | See “To migrate from disk-based fencing to server-based fencing using the installer” on page 61. |
| Migrating manually | See “To manually migrate from disk-based fencing to server-based fencing” on page 64. |

To migrate from disk-based fencing to server-based fencing using the installer

- 1 Make sure system-to-system communication is functioning properly.
- 2 Make sure that the SF Oracle RAC cluster is online and uses disk-based fencing.

```
# vxfenadm -d
```

For example, if SF Oracle RAC cluster uses disk-based fencing:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: SCSI3
Fencing SCSI3 Disk Policy: dmp
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

- 3 On any node in the cluster, start the `installsfrac` program with the `-fencing` option.

```
# /opt/VRTS/install/installsfrac<version> -fencing
```

Where `<version>` is the specific release version.

The `installsfrac` program starts with a copyright message and verifies the cluster information.

Note the location of log files which you can access in the event of any problem with the configuration process.

- 4 Confirm that you want to proceed with the I/O fencing configuration.
The installer verifies whether I/O fencing is configured in enabled mode.
- 5 Confirm that you want to reconfigure I/O fencing.
- 6 Review the I/O fencing configuration options that the program presents.
Type **4** to migrate to server-based I/O fencing.

```
Select the fencing mechanism to be configured in this  
Application Cluster [1-4,q] 4
```

- 7 From the list of coordination points that the installer presents, select the coordination points that you want to replace.

For example:

```
Select the coordination points you would like to remove  
from the currently configured coordination points:
```

- 1) emc_clariion0_62
- 2) emc_clariion0_65
- 3) emc_clariion0_66
- 4) All
- 5) None
- b) Back to previous menu

```
Enter the options separated by spaces: [1-5,b,q,?] (5)? 1 2
```

If you want to migrate to server-based fencing with no coordinator disks, type **4** to remove all the coordinator disks.

- 8 Enter the total number of new coordination points.
If you want to migrate to server-based fencing configuration with a mix of coordination points, the number you enter at this prompt must be a total of both the new CP servers and the new coordinator disks.
- 9 Enter the total number of new coordinator disks.
If you want to migrate to server-based fencing with no coordinator disks, type **0** at this prompt.
- 10 Enter the total number of virtual IP addresses or host names of the virtual IP address for each of the CP servers.
- 11 Enter the virtual IP addresses or host names of the virtual IP address for each of the CP servers.
- 12 Verify and confirm the coordination points information for the fencing reconfiguration.
- 13 Review the output as the installer performs the following tasks:

- Removes the coordinator disks from the coordinator disk group.
- Updates the application cluster details on each of the new CP servers.
- Prepares the `vxfenmode.test` file on all nodes.
- Runs the `vxfenswap` script.
 Note the location of the `vxfenswap.log` file which you can access in the event of any problem with the configuration process.
- Completes the I/O fencing migration.

14 If you want to send this installation information to Symantec, answer **y** at the prompt.

```
Would you like to send the information about this installation
to Symantec to help improve installation in the future? [y,n,q,?] (y) y
```

15 After the migration is complete, verify the change in the fencing mode.

```
# vxfenadm -d
```

For example, after the migration from disk-based fencing to server-based fencing in the SF Oracle RAC cluster:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: Customized
Fencing Mechanism: cps
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

16 Verify the current coordination points that the `vxfen` driver uses.

```
# vxfenconfig -l
```

To manually migrate from disk-based fencing to server-based fencing

- 1 Make sure system-to-system communication is functioning properly.
- 2 Make sure that the SF Oracle RAC cluster is online and uses disk-based fencing.

```
# vxfenadm -d
```

For example, if SF Oracle RAC cluster uses disk-based fencing:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: SCSI3
Fencing SCSI3 Disk Policy: dmp
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

- 3 Make sure that you performed the following tasks on the designated CP server:
 - Preparing to configure the new CP server.
 - Configuring the new CP server
 - Preparing the new CP server for use by the SF Oracle RAC cluster

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for the procedures.

- 4 Create a new `/etc/vxfenmode.test` file on each SF Oracle RAC cluster node with the fencing configuration changes such as the CP server information. Refer to the sample `vxfenmode` files in the `/etc/vxfen.d` folder.
- 5 From any node in the SF Oracle RAC cluster, start the `vx fenceswap` utility:

```
# vx fenceswap [-n]
```

- 6 Review the message that the utility displays and confirm whether you want to commit the change.
 - If you do not want to commit the new fencing configuration changes, press Enter or answer **n** at the prompt.

```
Do you wish to commit this change? [y/n] (default: n) n
```


The `vxfsenmode` utility rolls back the migration operation.

- If you want to commit the new fencing configuration changes, answer **y** at the prompt.

```
Do you wish to commit this change? [y/n] (default: n) y
```

If the utility successfully commits, the utility moves the `/etc/vxfsenmode.test` file to the `/etc/vxfsenmode` file.

7 After the migration is complete, verify the change in the fencing mode.

```
# vxfsenadm -d
```

For example, after the migration from disk-based fencing to server-based fencing in the SF Oracle RAC cluster:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: Customized
Fencing Mechanism: cps
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

8 Verify the current coordination points that the `vxfsen` driver uses.

```
# vxfsenconfig -l
```

Migrating from server-based to disk-based fencing in an online cluster

You can either use the installer or manually migrate from server-based fencing to disk-based fencing without incurring application downtime in the SF Oracle RAC clusters.

See [“About migrating between disk-based and server-based fencing configurations”](#) on page 60.

You can also use response files to migrate between fencing configurations.

See [“Migrating between fencing configurations using response files”](#) on page 71.

Warning: The cluster might panic if any node leaves the cluster membership before the coordination points migration operation completes.

This section covers the following procedures:

Migrating using the script-based installer	See “To migrate from server-based fencing to disk-based fencing using the installer” on page 66.
Migrating manually	See “To manually migrate from server-based fencing to disk-based fencing” on page 70.

To migrate from server-based fencing to disk-based fencing using the installer

- 1 Make sure system-to-system communication is functioning properly.
- 2 Make sure that the SF Oracle RAC cluster is configured to use server-based fencing.

```
# vxfenadm -d
```

For example, if the SF Oracle RAC cluster uses server-based fencing, the output appears similar to the following:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: Customized
Fencing Mechanism: cps
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

- 3** On any node in the cluster, start the `installsfrac` program with the `-fencing` option.

```
# /opt/VRTS/install/installsfrac<version> -fencing
```

Where `<version>` is the specific release version.

The `installsfrac` program starts with a copyright message and verifies the cluster information.

Note the location of log files which you can access in the event of any problem with the configuration process.

- 4** Confirm that you want to proceed with the I/O fencing configuration.
The installer verifies whether I/O fencing is configured in enabled mode.
- 5** Confirm that you want to reconfigure I/O fencing.
- 6** Review the I/O fencing configuration options that the program presents.
Type **4** to migrate to disk-based I/O fencing.

```
Select the fencing mechanism to be configured in this
Application Cluster [1-4,q] 4
```

- 7** From the list of coordination points that the installer presents, select the coordination points that you want to replace.

For example:

```
Select the coordination points you would like to remove
from the currently configured coordination points:
```

- 1) `emc_clariion0_62`
- 2) `[10.209.80.197]:14250, [10.209.80.199]:14300`
- 3) `[10.209.80.198]:14250`
- 4) All
- 5) None
- b) Back to previous menu

```
Enter the options separated by spaces: [1-5,b,q,?] (5)? 2 3
```

- 8** Enter the total number of new coordination points.
- 9** Enter the total number of new coordinator disks.

- 10 From the list of available disks that the installer presents, select two disks which you want to configure as coordinator disks.

For example:

List of available disks:

- 1) emc_clariion0_61
- 2) emc_clariion0_65
- 3) emc_clariion0_66
- b) Back to previous menu

Select 2 disk(s) as coordination points. Enter the disk options separated by spaces: [1-3,b,q]**2 3**

- 11 Verify and confirm the coordination points information for the fencing reconfiguration.

- 12 To migrate to disk-based fencing, select the I/O fencing mode as SCSI3.

Select the vxfen mode: [1-2,b,q,?] (1) **1**

The installer initializes the coordinator disks and the coordinator disk group, and depots the disk group. Press **Enter** to continue.

- 13 Review the output as the installer prepares the `vxfenmode.test` file on all nodes and runs the `vxfenswap` script.

Note the location of the `vxfenswap.log` file which you can access in the event of any problem with the configuration process.

The installer cleans up the application cluster information from the CP servers.

- 14 If you want to send this installation information to Symantec, answer **y** at the prompt.

Would you like to send the information about this installation to Symantec to help improve installation in the future? [y,n,q,?] (y) **y**

- 15** After the migration is complete, verify the change in the fencing mode.

```
# vxfenadm -d
```

For example, after the migration from server-based fencing to disk-based fencing in the SF Oracle RAC cluster:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: SCSI3
Fencing SCSI3 Disk Policy: dmp
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

- 16** Verify the current coordination points that the vxfen driver uses.

```
# vxfenconfig -l
```

To manually migrate from server-based fencing to disk-based fencing

- 1 Make sure system-to-system communication is functioning properly.
- 2 Make sure that the SF Oracle RAC cluster is online and uses server-based fencing.

```
# vxfenadm -d
```

For example, if SF Oracle RAC cluster uses server-based fencing:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: Customized
Fencing Mechanism: cps
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

- 3 Make sure that you performed the following preparatory tasks to configure disk-based fencing:

- Identifying disks to use as coordinator disks
- Setting up coordinator disk group
- Creating I/O configuration files

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for the procedures.

- 4 Create a new `/etc/vxfenmode.test` file with the fencing configuration changes such as the `scsi3` disk policy information.

Refer to the sample `vxfenmode` files in the `/etc/vxfen.d` folder.

- 5 From any node in the SF Oracle RAC cluster, start the `vxfenswap` utility:

```
# vxfenswap -g diskgroup [-n]
```

- 6 Review the message that the utility displays and confirm whether you want to commit the change.

- If you do not want to commit the new fencing configuration changes, press Enter or answer `n` at the prompt.

Do you wish to commit this change? [y/n] (default: n) **n**

The vxfenswap utility rolls back the migration operation.

- If you want to commit the new fencing configuration changes, answer **y** at the prompt.

Do you wish to commit this change? [y/n] (default: n) **y**

If the utility successfully commits, the utility moves the
 /etc/vxfenmode.test file to the /etc/vxfenmode file.

7 After the migration is complete, verify the change in the fencing mode.

```
# vxfenadm -d
```

For example, after the migration from server-based fencing to disk-based fencing in the SF Oracle RAC cluster:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: SCSI3
Fencing SCSI3 Disk Policy: dmp
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

8 Verify the current coordination points that the vxfen driver uses.

```
# vxfenconfig -l
```

Migrating between fencing configurations using response files

Typically, you can use the response file that the installer generates after you migrate between I/O fencing configurations. Edit these response files to perform an automated fencing reconfiguration in the SF Oracle RAC cluster.

To configure I/O fencing using response files

- 1 Make sure that SF Oracle RAC is configured.
- 2 Make sure system-to-system communication is functioning properly.

- 3 Make sure that the SF Oracle RAC cluster is online and uses either disk-based or server-based fencing.

```
# vxfenadm -d
```

For example, if SF Oracle RAC cluster uses disk-based fencing:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: SCSI3
Fencing SCSI3 Disk Policy: dmp
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

For example, if the SF Oracle RAC cluster uses server-based fencing:

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: Customized
Fencing Mechanism: cps
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```


- 4 Copy the response file to one of the cluster systems where you want to configure I/O fencing.

Review the sample files to reconfigure I/O fencing.

See [“Sample response file to migrate from disk-based to server-based fencing”](#) on page 73.

See [“Sample response file to migrate from server-based fencing to disk-based fencing”](#) on page 74.

See [“Sample response file to migrate from single CP server-based fencing to server-based fencing”](#) on page 74.

- 5 Edit the values of the response file variables as necessary.

See [“Response file variables to migrate between fencing configurations”](#) on page 74.

- 6 Start the I/O fencing reconfiguration from the system to which you copied the response file. For example:

```
# /opt/VRTS/install/installsffrac<version> -responsefile /tmp/  
\ response_file
```

Where *<version>* is the specific release version, and */tmp/response_file* is the response file's full path name.

Sample response file to migrate from disk-based to server-based fencing

The following is a sample response file to migrate from disk-based fencing with three coordinator disks to server-based fencing with one CP server and two coordinator disks:

```
$CFG{disks_to_remove}=[ qw(emc_clariion0_62) ];  
$CFG{fencing_cps}=[ qw(10.198.89.251) ];  
$CFG{fencing_cps_ports}{"10.198.89.204"}=14250;  
$CFG{fencing_cps_ports}{"10.198.89.251"}=14250;  
$CFG{fencing_cps_vips}{"10.198.89.251"}=[ qw(10.198.89.251 10.198.89.204) ]  
$CFG{fencing_ncp}=1;  
$CFG{fencing_option}=4;  
$CFG{opt}{configure}=1;  
$CFG{opt}{fencing}=1;  
$CFG{prod}="SFRAC60";  
$CFG{systems}=[ qw(sys1 sys2) ];  
$CFG{vcs_clusterid}=22462;  
$CFG{vcs_clustername}="clus1";
```

Sample response file to migrate from server-based fencing to disk-based fencing

The following is a sample response file to migrate from server-based fencing with one CP server and two coordinator disks to disk-based fencing with three coordinator disks:

```
$CFG{fencing_disks}=[ qw(emc_clariion0_66) ];  
$CFG{fencing_mode}="scsi3";  
$CFG{fencing_ncp}=1;  
$CFG{fencing_ndisks}=1;  
$CFG{fencing_option}=4;  
$CFG{opt}{configure}=1;  
$CFG{opt}{fencing}=1;  
$CFG{prod}="SFRAC60";  
$CFG{servers_to_remove}=[ qw([10.198.89.251]:14250) ];  
$CFG{systems}=[ qw(sys1 sys2) ];  
$CFG{vcs_clusterid}=42076;  
$CFG{vcs_clustername}="clus1";
```

Sample response file to migrate from single CP server-based fencing to server-based fencing

The following is a sample response file to migrate from single CP server-based fencing to server-based fencing with one CP server and two coordinator disks:

```
$CFG{fencing_disks}=[ qw(emc_clariion0_62 emc_clariion0_65) ];  
$CFG{fencing_dgname}="fendg";  
$CFG{fencing_scsi3_disk_policy}="dmp";  
$CFG{fencing_ncp}=2;  
$CFG{fencing_ndisks}=2;  
$CFG{fencing_option}=4;  
$CFG{opt}{configure}=1;  
$CFG{opt}{fencing}=1;  
$CFG{prod}="SFRAC60";  
$CFG{systems}=[ qw(sys1 sys2) ];  
$CFG{vcs_clusterid}=42076;  
$CFG{vcs_clustername}="clus1";
```

Response file variables to migrate between fencing configurations

Table 1-5 lists the response file variables that specify the required information to migrate between fencing configurations for SF Oracle RAC.

Table 1-5 Response file variables specific to migrate between fencing configurations

Variable	List or Scalar	Description
CFG {fencing_option}	Scalar	<p>Specifies the I/O fencing configuration mode.</p> <ul style="list-style-type: none"> ■ 1—Coordination Point Server-based I/O fencing ■ 2—Coordinator disk-based I/O fencing ■ 3—Disabled mode ■ 4—Fencing migration when the cluster is online <p>(Required)</p>
CFG {fencing_reusedisk}	Scalar	<p>If you migrate to disk-based fencing or to server-based fencing that uses coordinator disks, specifies whether to use free disks or disks that already belong to a disk group.</p> <ul style="list-style-type: none"> ■ 0—Use free disks as coordinator disks ■ 1—Use disks that already belong to a disk group as coordinator disks (before configuring these as coordinator disks, installer removes the disks from the disk group that the disks belonged to.) <p>(Required if your fencing configuration uses coordinator disks)</p>
CFG {fencing_ncp}	Scalar	<p>Specifies the number of new coordination points to be added.</p> <p>(Required)</p>
CFG {fencing_ndisks}	Scalar	<p>Specifies the number of disks in the coordination points to be added.</p> <p>(Required if your fencing configuration uses coordinator disks)</p>

Table 1-5 Response file variables specific to migrate between fencing configurations (*continued*)

Variable	List or Scalar	Description
CFG {fencing_disks}	List	Specifies the disks in the coordination points to be added. (Required if your fencing configuration uses coordinator disks)
CFG {fencing_dgname}	Scalar	Specifies the disk group that the coordinator disks are in. (Required if your fencing configuration uses coordinator disks)
CFG {fencing_scsi3_disk_policy}	Scalar	Specifies the disk policy that the disks must use. (Required if your fencing configuration uses coordinator disks)
CFG {fencing_cps}	List	Specifies the CP servers in the coordination points to be added. (Required for server-based fencing)
CFG {fencing_cps_vips}{vip1}	List	Specifies the virtual IP addresses or the fully qualified host names of the new CP server. (Required for server-based fencing)
CFG {fencing_cps_ports}{vip}	Scalar	Specifies the port that the virtual IP of the new CP server must listen on. If you do not specify, the default value is 14250. (Optional)
CFG {servers_to_remove}	List	Specifies the CP servers in the coordination points to be removed.
CFG {disks_to_remove}	List	Specifies the disks in the coordination points to be removed

About making CP server highly available

If you want to configure a multi-node CP server cluster, install and configure SFHA on the CP server nodes. Otherwise, install and configure VCS on the single node.

In both the configurations, VCS provides local start and stop of the CP server process, taking care of dependencies such as NIC, IP address, and so on. Moreover, VCS also serves to restart the CP server process in case the process faults.

VCS can use multiple network paths to access a CP server. If a network path fails, CP server does not require a restart and continues to listen on one of the other available virtual IP addresses.

To make the CP server process highly available, you must perform the following tasks:

- Install and configure SFHA on the CP server systems.
- Configure the CP server process as a failover service group.
- Configure disk-based I/O fencing to protect the shared CP server database.

Note: Symantec recommends that you do not run any other applications on the single node or SFHA cluster that is used to host CP server.

A single CP server can serve multiple SF Oracle RAC clusters. A common set of CP servers can serve up to 128 SF Oracle RAC clusters.

About secure communication between the SF Oracle RAC cluster and CP server

In a data center, TCP/IP communication between the SF Oracle RAC cluster (application cluster) and CP server must be made secure. The security of the communication channel involves encryption, authentication, and authorization.

The CP server node or cluster needs to confirm the authenticity of the SF Oracle RAC cluster nodes that communicate with it as a coordination point and only accept requests from known SF Oracle RAC cluster nodes. Requests from unknown clients are rejected as non-authenticated. Similarly, the fencing framework in SF Oracle RAC cluster must confirm that authentic users are conducting fencing operations with the CP server.

Entities on behalf of which authentication is done, are referred to as principals. On the SF Oracle RAC cluster nodes, the current VCS installer creates the Authentication Server credentials on each node in the cluster. It also creates vcsauthserver which authenticates the credentials. The installer then proceeds to start VCS in secure mode.

Typically, in an existing VCS cluster with security configured, vcsauthserver runs on each cluster node.

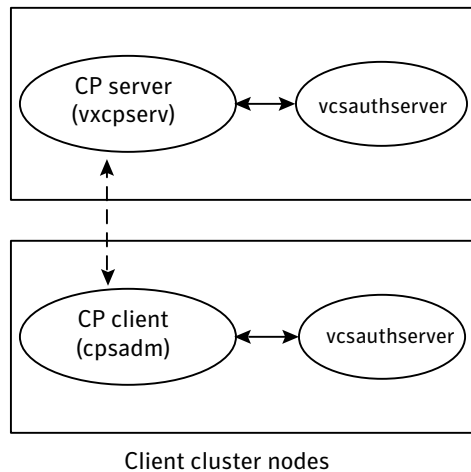
How secure communication between the CP servers and the SF Oracle RAC clusters work

CP server and SF Oracle RAC cluster (application cluster) node communication involve the following entities:

- vxcpserv for the CP server
- cpsadm for the SF Oracle RAC cluster node

Figure 1-12 displays a schematic of the end-to-end communication flow with security enabled on CP server and SF Oracle RAC clusters (application clusters).

Figure 1-12 End-To-end communication flow with security enabled on CP server and SF Oracle RAC clusters



Communication flow between CP server and SF Oracle RAC cluster nodes with security configured on them is as follows:

- Initial setup:
Identities of CP server and SF Oracle RAC cluster nodes are configured on respective nodes by the VCS installer.

Note: At the time of fencing configuration, the installer establishes trust between the CP server and the application cluster so that vxcpserv process can authenticate requests from the application cluster nodes. If you manually configured I/O fencing, then you must set up trust between the CP server and the application cluster.

The `cpsadm` command gets the user name, domain type from the environment variables `CPS_USERNAME`, `CPS_DOMAINTYPE`. The user is expected to export these variables before running the `cpsadm` command manually. The customized fencing framework exports these environment variables internally before running the `cpsadm` commands.

The CP server process (`vxcpserv`) uses its own user (`CPSERVER`) which is added to the local `vcsauthserver`.

- Getting credentials from authentication broker:
The `cpsadm` command tries to get the existing credentials that are present on the local node. The installer generates these credentials during fencing configuration.
The `vxcpserv` process tries to get the existing credentials that are present on the local node. The installer generates these credentials when it enables security.
- Communication between CP server and SF Oracle RAC cluster nodes:
After the CP server establishes its credential and is up, it becomes ready to receive data from the clients. After the `cpsadm` command obtains its credentials and authenticates CP server credentials, `cpsadm` connects to the CP server. Data is passed over to the CP server.
- Validation:
On receiving data from a particular SF Oracle RAC cluster node, `vxcpserv` validates its credentials. If validation fails, then the connection request data is rejected.

Security configuration details on CP server and SF Oracle RAC cluster

This section discusses the security configuration details for the CP server and SF Oracle RAC cluster (application cluster).

Settings in secure mode

The following are the settings for secure communication between the CP server and SF Oracle RAC cluster:

- CP server settings:
Installer creates a user with the following values:
 - username: `CPSERVER`
 - domainname: `VCS_SERVICES@cluster_uuid`
 - domaintype: `vx`

Run the following commands on the CP server to verify the settings:

```
# export EAT_DATA_DIR=/var/VRTSvcs/vcsauth/data/CPSERVER
```

```
# /opt/VRTScps/bin/cpsat showcred
```

Note: The CP server configuration file (`/etc/vxcps.conf`) must not contain a line specifying **security=0**. If there is no line specifying `security` parameter or if there is a line specifying **security=1**, CP server with security is enabled (which is the default).

■ SF Oracle RAC cluster node(s) settings:

On SF Oracle RAC cluster, the installer creates a user for `cpsadm` during fencing configuration with the following values:

- username: CPSADM
- domainname: VCS_SERVICES@*cluster_uuid*
- domaintype: vx

Run the following commands on the SF Oracle RAC cluster node(s) to verify the security settings:

```
# export EAT_DATA_DIR=/var/VRTSvc/vcsauth/data/CPSADM  
  
# /opt/VRTScps/bin/cpsat showcred
```

The users described above are used only for authentication for the communication between the CP server and the SF Oracle RAC cluster nodes.

For CP server's authorization, customized fencing framework on the SF Oracle RAC cluster uses the following user if security is configured:

CPSADM@VCS_SERVICES@*cluster_uuid*

where *cluster_uuid* is the application cluster's universal unique identifier.

For each SF Oracle RAC cluster node, this user must be registered on the CP server database before fencing starts on the SF Oracle RAC cluster node(s). This can be verified by issuing the following command:

```
# cpsadm -s cp_server -a list_users
```

The following is an example of the command output:

```
Username/Domain Type  
CPSADM@VCS_SERVICES@77a2549c-1dd2-11b2-88d6-00306e4b2e0b/vx  
  
Cluster Name / UUID                               Role  
cluster1/{77a2549c-1dd2-11b2-88d6-00306e4b2e0b} Operator
```

Note: The configuration file (`/etc/vxfenmode`) on each client node must not contain a line specifying **security=0**. If there is no line specifying `security` parameter or if there is a line specifying **security=1**, client node starts with security enabled (which is the default).

Settings in non-secure mode

In non-secure mode, only authorization is provided on the CP server. Passwords are not requested. Authentication and encryption are not provided. User credentials of “`cpsclient@hostname`” of “`vx`” domain type are used by the customized fencing framework for communication between CP server or SF Oracle RAC cluster node(s).

For each SF Oracle RAC cluster node, this user must be added on the CP server database before fencing starts on the SF Oracle RAC cluster node(s). The user can be verified by issuing the following command:

```
# cpsadm -s cpserver -a list_users
```

The following is an example of the command output:

Username/Domain	Type	Cluster Name / UUID	Role
<code>cpsclient@galaxy/vx</code>		<code>cluster1 / {f0735332-e3709c1c73b9}</code>	Operator

Note: In non-secure mode, CP server configuration file (`/etc/vxcps.conf`) should contain a line specifying **security=0**. Similarly, on each SF Oracle RAC cluster node the configuration file (`/etc/vxfenmode`) should contain a line specifying **security=0**.

Oracle RAC components

This section provides a brief description of Oracle Clusterware/Grid Infrastructure, the Oracle Cluster Registry, application resources, and the voting disk.

Note: Refer to the Oracle RAC documentation for additional information.

Oracle Clusterware/Grid Infrastructure

Oracle Clusterware/Grid Infrastructure manages Oracle cluster-related functions including membership, group services, global resource management, and databases. Oracle Clusterware/Grid Infrastructure is required for every Oracle RAC instance.

Oracle Clusterware/Grid Infrastructure requires the following major components:

- A cluster interconnect that allows for cluster communications
- A private virtual IP address for cluster communications over the interconnect
- A public virtual IP address for client connections
- For Oracle 11g Release 2, a public IP address as a Single Client Access Name (SCAN) address on the Domain Name Server (DNS) for round robin resolution to three IP addresses (recommended) or at least one IP address
- Shared storage accessible by each node

Co-existence with VCS

Oracle Clusterware/Grid Infrastructure supports co-existence with vendor clusterwares such as Veritas Cluster Server. When you install Oracle Clusterware/Grid Infrastructure on an SF Oracle RAC cluster, Oracle Clusterware/Grid Infrastructure detects the presence of VCS by checking the presence of the Veritas membership module (VCSMM) library. It obtains the list of nodes in the cluster from the VCSMM library at the time of installation.

When a node fails to respond across the interconnect, Oracle Clusterware/Grid Infrastructure waits before evicting another node from the cluster. This wait-time is defined by the CSS miss-count value. Oracle Clusterware/Grid Infrastructure sets the CSS miss-count parameter to a larger value (600 seconds) in the presence of VCS. This value is much higher than the LLT peer inactivity timeout interval. Thus, in the event of a network split-brain, the two clusterwares, VCS and Oracle Clusterware/Grid Infrastructure, do not interfere with each other's decisions on which nodes remain in the cluster. Veritas I/O fencing is allowed to decide on the surviving nodes first, followed by Oracle Clusterware/Grid Infrastructure.

VCS uses LLT for communicating between cluster nodes over private interconnects while Oracle Clusterware/Grid Infrastructure uses private IP addresses configured over the private interconnects to communicate between the cluster nodes. To coordinate membership changes between VCS and Oracle Clusterware/Grid Infrastructure, it is important to configure the Oracle Clusterware/Grid Infrastructure private IP address over the network interfaces used by LLT. VCS uses the CSSD agent to start, stop, and monitor Oracle Clusterware/Grid Infrastructure. The CSSD agent ensures that the OCR, the voting disk, and the private IP address resources required by Oracle Clusterware/Grid Infrastructure are brought online by VCS before Oracle Clusterware/Grid Infrastructure starts. This prevents the premature startup of Oracle Clusterware/Grid Infrastructure, which causes cluster failures.

Oracle Cluster Registry

The Oracle Cluster Registry (OCR) contains cluster and database configuration and state information for Oracle RAC and Oracle Clusterware/Grid Infrastructure.

The information maintained in the OCR includes:

- The list of nodes
- The mapping of database instances to nodes
- Oracle Clusterware application resource profiles
- Resource profiles that define the properties of resources under Oracle Clusterware/Grid Infrastructure control
- Rules that define dependencies between the Oracle Clusterware/Grid Infrastructure resources
- The current state of the cluster

For versions prior to Oracle RAC 11g Release 2, the OCR data exists on a shared raw volume or a cluster file system that is accessible to each node.

In Oracle RAC 11g Release 2, the OCR data exists on ASM or a cluster file system that is accessible to each node.

Use CVM mirrored volumes to protect OCR data from failures. Oracle Clusterware/Grid Infrastructure faults nodes if OCR is not accessible because of corruption or disk failure. Oracle automatically backs up OCR data. You can also export the OCR contents before making configuration changes in Oracle Clusterware/Grid Infrastructure. This way, if you encounter configuration problems and are unable to restart Oracle Clusterware/Grid Infrastructure, you can restore the original contents.

Consult the Oracle documentation for instructions on exporting and restoring OCR contents.

Application resources

Oracle Clusterware/Grid Infrastructure application resources are similar to VCS resources. Each component controlled by Oracle Clusterware/Grid Infrastructure is defined by an application resource, including databases, instances, services, and node applications.

Unlike VCS, Oracle Clusterware/Grid Infrastructure uses separate resources for components that run in parallel on multiple nodes.

Resource profiles

Resources are defined by profiles, which are similar to the attributes that define VCS resources. The OCR contains application resource profiles, dependencies, and status information.

Oracle Clusterware/Grid Infrastructure node applications

Oracle Clusterware/Grid Infrastructure uses these node application resources to manage Oracle components, such as databases, listeners, and virtual IP addresses. Node application resources are created during Oracle Clusterware/Grid Infrastructure installation.

Voting disk

The voting disk is a heartbeat mechanism used by Oracle Clusterware/Grid Infrastructure to maintain cluster node membership.

In versions prior to Oracle RAC 11g Release 2, voting disk data exists on a shared raw volume or a cluster file system that is accessible to each node.

In Oracle RAC 11g Release 2, voting disk data exists on ASM or on a cluster file system that is accessible to each node.

The Oracle Clusterware Cluster Synchronization Service daemon (ocssd) provides cluster node membership and group membership information to RAC instances. On each node, ocssd processes write a heartbeat to the voting disk every second. If a node is unable to access the voting disk, Oracle Clusterware/Grid Infrastructure determines the cluster is in a split-brain condition and panics the node in a way that only one sub-cluster remains.

Oracle Disk Manager

SF Oracle RAC requires Oracle Disk Manager (ODM), a standard API published by Oracle for support of database I/O. SF Oracle RAC provides a library for Oracle to use as its I/O library.

ODM architecture

When the Veritas ODM library is linked, Oracle is able to bypass all caching and locks at the file system layer and to communicate directly with raw volumes. The SF Oracle RAC implementation of ODM generates performance equivalent to performance with raw devices while the storage uses easy-to-manage file systems.

All ODM features can operate in a cluster environment. Nodes communicate with each other before performing any operation that could potentially affect another node. For example, before creating a new data file with a specific name, ODM checks with other nodes to see if the file name is already in use.

Veritas ODM performance enhancements

Veritas ODM enables the following performance benefits provided by Oracle Disk Manager:

- Locking for data integrity.
- Few system calls and context switches.
- Increased I/O parallelism.
- Efficient file creation and disk allocation.

Databases using file systems typically incur additional overhead:

- Extra CPU and memory usage to read data from underlying disks to the file system cache. This scenario requires copying data from the file system cache to the Oracle cache.
- File locking that allows for only a single writer at a time. Allowing Oracle to perform locking allows for finer granularity of locking at the row level.
- File systems generally go through a standard Sync I/O library when performing I/O. Oracle can make use of Kernel Async I/O libraries (KAIO) with raw devices to improve performance.

ODM communication

ODM uses port d to communicate with ODM instances on other nodes in the cluster.

RAC extensions

Oracle RAC relies on several support services provided by VCS. Key features include Veritas Cluster Server Membership Manager (VCSMM) and Veritas Cluster Server Inter-Process Communication (VCSIPC), and LLT Multiplexer (LMX).

Veritas Cluster Server membership manager

To protect data integrity by coordinating locking between RAC instances, Oracle must know which instances actively access a database. Oracle provides an API called `skgxn` (system kernel generic interface node membership) to obtain information on membership. SF Oracle RAC implements this API as a library linked to Oracle after you install Oracle RAC. Oracle uses the linked `skgxn` library to make `ioctl` calls to VCSMM, which in turn obtains membership information for clusters and instances by communicating with GAB on port o.

LLT multiplexer

The LMX module is a kernel module designed to receive communications from the `skgxp` module and pass them on to the correct process on the correct instance on other nodes. The LMX module "multiplexes" communications between multiple processes on other nodes. LMX uses all LLT links to send VCS IPC (Oracle Cache Fusion) data. LMX leverages all features of LLT, including load balancing and fault resilience.

Note: The LLT multiplexer (LMX) is not supported with Oracle RAC 11g.

Veritas Cluster Server inter-process communication

To coordinate access to a single database by multiple instances, Oracle uses extensive communications between nodes and instances. Oracle uses Inter-Process Communications (VCSIPC) for Global Enqueue Service locking traffic and Global Cache Service cache fusion. SF Oracle RAC uses LLT to support VCSIPC in a cluster and leverages its high-performance and fault-resilient capabilities.

Oracle has an API for VCSIPC, System Kernel Generic Interface Inter-Process Communications (`skgxp`), that isolates Oracle from the underlying transport mechanism. As Oracle conducts communication between processes, it does not need to know how data moves between systems; the cluster implementer can create the highest performance for internode communications without Oracle reconfiguration.

SF Oracle RAC provides a library to implement the `skgxp` functionality. This module communicates with the LLT Multiplexer (LMX) using `ioctl` calls.

Oracle and cache fusion traffic

Private IP addresses are required by Oracle for cache fusion traffic.

Depending on the version of Oracle RAC you want to install, you have the following options for setting up your private network configuration:

- | | |
|----------------|--|
| Oracle RAC 10g | Use either Oracle UDP IPC or VCSIPC/LMX/LLT for the database cache fusion traffic.

By default, the database cache fusion traffic is configured to use VCSIPC/LMX/LLT. |
| Oracle RAC 11g | You must use UDP IPC for the database cache fusion traffic. |

Periodic health evaluation of SF Oracle RAC clusters

SF Oracle RAC provides a health check utility that evaluates the components and configuration in an SF Oracle RAC cluster. The utility when invoked gathers real-time operational information on cluster components and displays the report on your system console. You must run the utility on each node in the cluster.

The utility evaluates the health of the following components:

- Low Latency Transport (LLT)
- LLT Multiplexer (LMX)
- VCSMM
- I/O fencing
- Oracle Clusterware/Grid Infrastructure
- PrivNIC and MultiPrivNIC

The health check utility is installed at

`/opt/VRTSvcs/rac/healthcheck/healthcheck` during the installation of SF Oracle RAC.

The utility determines the health of the components by gathering information from configuration files or by reading the threshold values set in the health check configuration file `/opt/VRTSvcs/rac/healthcheck/healthcheck.cf`

The utility displays a warning message when the operational state of the component violates the configured settings for the component or when the health score approaches or exceeds the threshold value. You can modify the health check configuration file to set the threshold values to the desired level.

The health checks for the VCSMM, LMX, I/O fencing, PrivNIC, MultiPrivNIC, and Oracle Clusterware components do not use threshold settings.

Note: You must set the `ORACLE_HOME` and `CRS_HOME` parameters in the configuration file as appropriate for your setup.

[Table 1-6](#) provides guidelines on changing the threshold values.

Table 1-6 Setting threshold values

Requirement	Setting the threshold
To detect warnings early	Reduce the corresponding threshold value.

Table 1-6 Setting threshold values (*continued*)

Requirement	Setting the threshold
To suppress warnings	Increase the corresponding threshold value. Caution: Using very high threshold values can prevent the health check utility from forecasting potential problems in the cluster. Exercise caution with high values.

Note: You can schedule periodic health evaluation of your clusters, by scheduling the utility to run as a cron job.

See [“Scheduling periodic health checks for your SF Oracle RAC cluster”](#) on page 109.

For detailed information on the list of health checks performed for each component, see the appendix *List of SF Oracle RAC health checks*.

About Virtual Business Services

Virtual Business Services provide continuous high availability and reduce frequency and duration of service disruptions for multi-tier business applications running on heterogeneous operating systems and virtualization technologies. A Virtual Business Service represents the multi-tier application as a single consolidated entity and builds on the high availability and disaster recovery provided for the individual tiers by Symantec products such as Veritas Cluster Server and Symantec ApplicationHA. Additionally, a Virtual Business Service can also represent all the assets used by the service such as arrays, hosts, and file systems, though they are not migrated between server tiers. A Virtual Business Service provides a single consolidated entity that represents a multi-tier business service in its entirety. Application components that are managed by Veritas Cluster Server or Symantec ApplicationHA can be actively managed through a Virtual Business Service.

You can configure and manage Virtual Business Services created in Veritas Operations Manager by using Veritas Operations Manager Virtual Business Services Availability Add-on. Besides providing all the functionality that was earlier available through Business Entity Operations Add-on, VBS Availability Add-on provides the additional ability to configure fault dependencies between the components of the multi-tier application.

Note: All the Application Entities that were created using Veritas Operations Manager Business Entity Operations Add-on versions 3.1 and 4.0 are available as Virtual Business Services after you deploy the VBS Availability Add-on in Veritas Operations Manager 5.0. Veritas Operations Manager 5.0 is a prerequisite for running Virtual Business Services.

Features of Virtual Business Services

You can use the VBS Availability Add-on to perform the following tasks:

- Start Virtual Business Services from the Veritas Operations Manager console. When a Virtual Business Service starts, its associated service groups are brought online.
- Stop Virtual Business Services from the Veritas Operations Manager console. When a Virtual Business Service stops, its associated service groups are taken offline.
Applications that are under the control of Symantec ApplicationHA can be part of a Virtual Business Service. Symantec ApplicationHA enables starting, stopping, and monitoring of an application within a virtual machine. If applications are hosted on VMware virtual machines, you can configure the virtual machines to automatically start or stop when you start or stop the Virtual Business Service.
- Establish service group relationships and set the order to bring service groups online and to take them offline. It ensures that the service groups from different clusters are brought online or taken offline in the correct order. This order is governed by the service group's relationships with other service groups, which are referred to as child service groups. Setting the correct order of service group dependency is critical to achieve business continuity and high availability.
- Establish service group relationships and specify the required reaction of an application component to a high availability event in an underlying tier.
- Manage the Virtual Business Service from Veritas Operations Manager or from the clusters participating in the Virtual Business Service.
- Recover the entire Virtual Business Service to a remote site when a disaster occurs.

However, the following operations cannot be managed using VBS Availability Add-on:

- The service group operations that are performed using the Veritas Cluster Server management console.

- The service group operations that are performed using the Veritas Cluster Server command-line interface.
- The service group operations that are performed using the Veritas Cluster Server Java console.
- VBS Availability Add-on is not supported for composite Virtual Business Services. You can use it only for Virtual Business Services.

Note: You must install the VRTSvbs depot on the cluster nodes to enable fault management and to administer the Virtual Business Service from the participating clusters.

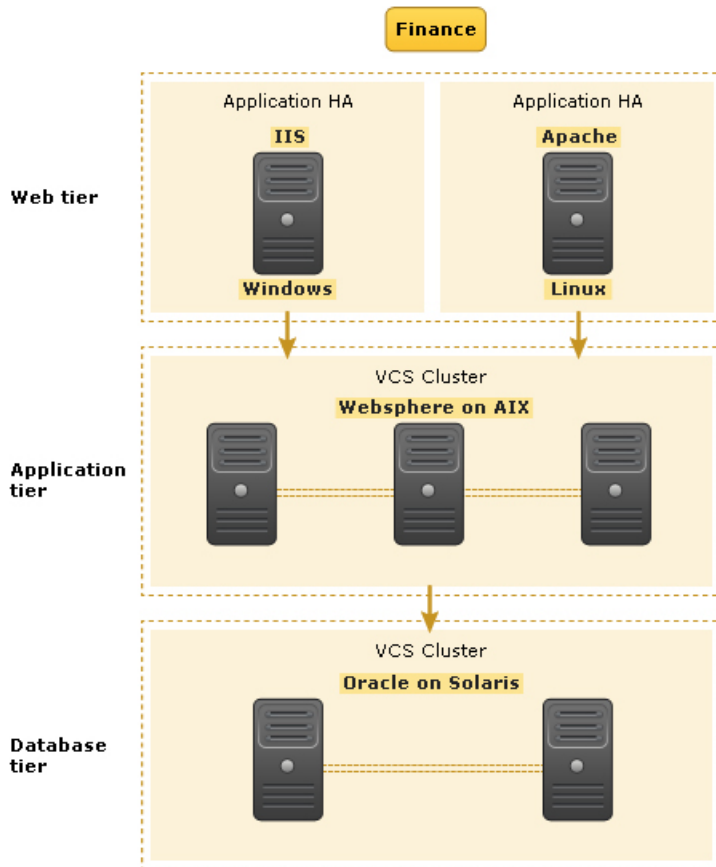
Sample Virtual Business Service configuration

This section provides a sample Virtual Business Service configuration comprising a multi-tier application. [Figure 1-13](#) shows a Finance application that is dependent on components that run on three different operating systems and on three different clusters.

- Databases such as Oracle running on Solaris operating systems form the database tier.
- Middleware applications such as WebSphere running on AIX operating systems form the middle tier.
- Web applications such as Apache and IIS running on Windows and Linux virtual machines form the Web tier. This tier is composed of ApplicationHA nodes.

Each tier can have its own high availability mechanism. For example, you can use Veritas Cluster Server for the databases and middleware applications, and Symantec ApplicationHA for the Web servers.

Figure 1-13 Sample Virtual Business Service configuration



Each time you start the Finance business application, typically you need to bring the components online in the following order – Oracle database, WebSphere, Apache and IIS. In addition, you must bring the virtual machines online before you start the Web tier. To stop the Finance application, you must take the components offline in the reverse order. From the business perspective, the Finance service is unavailable if any of the tiers becomes unavailable.

When you configure the Finance application as a Virtual Business Service, you can specify that the Oracle database must start first, followed by WebSphere and the Web servers. The reverse order automatically applies when you stop the Virtual Business Service. When you start or stop the Virtual Business Service, the components of the service are started or stopped in the defined order.

For more information about Virtual Business Services, refer to the *Virtual Business Service–Availability User’s Guide*.

See “[About Virtual Business Services](#)” on page 88.

About Veritas Operations Manager

Veritas Operations Manager provides a centralized management console for Veritas Storage Foundation and High Availability products. You can use Veritas Operations Manager to monitor, visualize, and manage storage resources and generate reports.

Symantec recommends using Veritas Operations Manager (VOM) to manage Storage Foundation and Cluster Server environments.

You can download Veritas Operations Manager at no charge at <http://go.symantec.com/vom>.

Refer to the Veritas Operations Manager documentation for installation, upgrade, and configuration instructions.

The Veritas Enterprise Administrator (VEA) console is no longer packaged with Storage Foundation products. If you want to continue using VEA, a software version is available for download from http://go.symantec.com/vcsm_download. Veritas Storage Foundation Management Server is deprecated.

If you want to manage a single cluster using Cluster Manager (Java Console), a version is available for download from http://go.symantec.com/vcsm_download. You cannot manage the new features of this release using the Java Console. Veritas Cluster Server Management Console is deprecated.

About Symantec Operations Readiness Tools

[Symantec Operations Readiness Tools \(SORT\)](#) is a Web site that automates and simplifies some of the most time-consuming administrative tasks. SORT helps you manage your datacenter more efficiently and get the most out of your Symantec products.

Among its broad set of features, SORT lets you do the following:

- Generate server-specific reports that describe how to prepare your servers for installation or upgrade of Symantec enterprise products.
- Access a single site with the latest production information, including patches, agents, and documentation.
- Create automatic email notifications for changes in patches, documentation, and array-specific modules.

To access SORT, go to:

<https://sort.symantec.com>

Administering SF Oracle RAC and its components

This chapter includes the following topics:

- [Administering SF Oracle RAC](#)
- [Administering VCS](#)
- [Administering I/O fencing](#)
- [Administering the CP server](#)
- [Administering CFS](#)
- [Administering CVM](#)
- [Administering SF Oracle RAC global clusters](#)

Administering SF Oracle RAC

This section provides instructions for the following SF Oracle RAC administration tasks:

- Setting the environment variables
See [“Setting the environment variables for SF Oracle RAC”](#) on page 96.
- Starting or stopping SF Oracle RAC on each node
See [“Starting or stopping SF Oracle RAC on each node”](#) on page 97.
- Applying Oracle patches on SF Oracle RAC nodes
See [“Applying Oracle patches on SF Oracle RAC nodes”](#) on page 102.
- Installing Veritas Volume Manager or Veritas File System patches or ODM patches on SF Oracle RAC nodes

See [“Installing Veritas Volume Manager, Veritas File System, or ODM patches on SF Oracle RAC nodes”](#) on page 103.

- Applying operating system updates on SF Oracle RAC nodes
See [“Applying operating system updates on SF Oracle RAC nodes”](#) on page 103.
- Determining the LMX traffic for each database in an SF Oracle RAC cluster
See [“Determining the LMX traffic for each database in an SF Oracle RAC cluster”](#) on page 104.
- Adding storage to an SF Oracle RAC cluster
See [“Adding storage to an SF Oracle RAC cluster”](#) on page 106.
- Recovering from storage failure
See [“Recovering from storage failure”](#) on page 107.
- Backing up and restoring the Oracle database using Symantec NetBackup
See [“Backing up and restoring Oracle database using Symantec NetBackup”](#) on page 107.
- Enhancing the performance of SF Oracle RAC clusters
See [“Enhancing the performance of SF Oracle RAC clusters”](#) on page 108.
- Creating snapshots for offhost processing
See [“Creating snapshots for offhost processing”](#) on page 109.
- Managing database storage efficiently using SmartTier
See [“Managing database storage efficiently using SmartTier”](#) on page 109.
- Optimizing database storage using Thin Provisioning and SmartMove
See [“Optimizing database storage using Thin Provisioning and SmartMove”](#) on page 109.
- Scheduling periodic health checks for your SF Oracle RAC cluster
See [“Scheduling periodic health checks for your SF Oracle RAC cluster”](#) on page 109.
- Verifying the nodes in an SF Oracle RAC cluster
See [“Verifying the nodes in an SF Oracle RAC cluster”](#) on page 110.

If you encounter issues while administering SF Oracle RAC, refer to the troubleshooting section for assistance.

See [“About troubleshooting SF Oracle RAC”](#) on page 179.

Setting the environment variables for SF Oracle RAC

Set the MANPATH variable in the .profile file (or other appropriate shell setup file for your system) to enable viewing of manual pages.

Based on the shell you use, type one of the following:

```
For sh, ksh, or bash    # MANPATH=/usr/share/man:/opt/VRTS/man
                        # export MANPATH
```

Set the PATH environment variable in the .profile file (or other appropriate shell setup file for your system) on each system to include installation and other commands.

Note: Do not define \$ORACLE_HOME/lib in LIBPATH for root user. You should define \$ORACLE_HOME/lib in LIBPATH for the oracle user.

Based on the shell you use, type one of the following:

```
For sh, ksh, or bash    # PATH=/usr/sbin:/sbin:/usr/bin:\
                        /opt/VRTS/bin:\
                        $PATH; export PATH
```

Starting or stopping SF Oracle RAC on each node

You can start or stop SF Oracle RAC on each node in the cluster using the SF Oracle RAC installer or manually.

To start SF Oracle RAC	Using installer: See “Starting SF Oracle RAC using the SF Oracle RAC installer” on page 98. Manual: See “Starting SF Oracle RAC manually on each node” on page 98.
To stop SF Oracle RAC	Using installer: See “Stopping SF Oracle RAC using the SF Oracle RAC installer” on page 99. Manual: See “Stopping SF Oracle RAC manually on each node” on page 99.

Starting SF Oracle RAC using the SF Oracle RAC installer

Run the SF Oracle RAC installer with the `-start` option to start SF Oracle RAC on each node.

Note: Start SF Oracle RAC on all nodes in the cluster. Specifying only some of the nodes in the cluster may cause some of the components that depend on GAB seeding to fail.

To start SF Oracle RAC using the SF Oracle RAC installer

- 1 Log into one of the nodes in the cluster as the root user.
- 2 Start SF Oracle RAC:

```
# /opt/VRTS/install/installsfrac<version> -start sys1 sys2
```

Where `<version>` is the specific release version.

Starting SF Oracle RAC manually on each node

Perform the steps in the following procedures to start SF Oracle RAC manually on each node.

To start SF Oracle RAC manually on each node

- 1 Log into each node as the root user.
- 2 Start LLT:

```
# /sbin/init.d/llt start
```

- 3 Start GAB:

```
# /sbin/init.d/gab start
```

- 4 Start fencing:

```
# /sbin/init.d/vxfen start
```

- 5 Start VCSMM:

```
# /sbin/init.d/vcsmm start
```

- 6 Start LMX:

```
# /sbin/init.d/lmx start
```

7 Start ODM:

```
# /sbin/init.d/odm start
```

8 Start VCS, CVM, and CFS:

```
# hstart
```

9 Verify that all GAB ports are up and running:

```
# gabconfig -a
```

```
GAB Port Memberships
```

```
=====
Port a gen 564004 membership 01
Port b gen 564008 membership 01
Port d gen 564009 membership 01
Port f gen 564024 membership 01
Port h gen 56401a membership 01
Port o gen 564007 membership 01
Port u gen 564021 membership 01
Port v gen 56401d membership 01
Port w gen 56401f membership 01
Port y gen 56401c membership 01
```

Stopping SF Oracle RAC using the SF Oracle RAC installer

Run the SF Oracle RAC installer with the `-stop` option to stop SF Oracle RAC on each node.

To stop SF Oracle RAC using the SF Oracle RAC installer

1 Log into one of the nodes in the cluster as the root user.

2 Stop VCS:

```
# hstop -all
```

3 Stop SF Oracle RAC:

```
# /opt/VRTS/install/installsfrac<version> -stop sys1 sys2
```

Stopping SF Oracle RAC manually on each node

Perform the steps in the following procedures to stop SF Oracle RAC manually on each node.

To stop SF Oracle RAC manually on each node

1 Stop the Oracle database.

If the Oracle RAC instance is managed by VCS, log in as the root user and take the service group offline:

```
# hagrps -offline oracle_group -sys node_name
```

If the Oracle database instance is not managed by VCS, log in as the Oracle user on one of the nodes and shut down the instance:

For Oracle RAC 11.2.0.2:

```
$ srvctl stop instance -d db_name \  
-n node_name
```

For Oracle RAC 11.2.0.1 and earlier versions:

```
$ srvctl stop instance -d db_name \  
-i instance_name
```

2 Stop all applications that are not configured under VCS but dependent on Oracle RAC or resources controlled by VCS. Use native application commands to stop the application.

3 Unmount the non-system CFS file systems that are not managed by VCS.

- Make sure that no processes are running which make use of mounted shared file system. To verify that no processes use the VxFS or CFS mount point:

```
# mount -v | grep vxfs | grep cluster
```

```
# fuser -cu /mount_point
```

- Unmount the non-system CFS file system:

```
# umount /mount_point
```

4 Take the VCS service groups offline:

```
# hagrps -offline group_name -sys node_name
```

Verify that the VCS service groups are offline:

```
# hagrps -state group_name
```

5 Unmount the non-system VxFS file systems that are not managed by VCS.

- Make sure that no processes are running which make use of mounted shared file system. To verify that no processes use the VxFS or CFS mount point:

```
# mount -v | grep vxfs
# fuser -cu /mount_point
```

- Unmount the non-system VxFS file system:

```
# umount /mount_point
```

- 6 If you created local VxFS mount points on VxVM volumes and added them to `/etc/fstab`, comment out the mount point entries in the file.
- 7 Verify that no VxVM volumes (other than VxVM boot volumes) remain open. Stop any open volumes that are not managed by VCS.
- 8 Unmount the VxFS file systems (except `/` and `/tmp`) that are not managed by VCS.

Make sure that no processes are running, which make use of mounted shared file system or shared volumes:

```
# mount -v | grep vxfs
# fuser -cu /mount_point
```

- 9 Stop VCS, CVM and CFS:

```
# hastop -local
```

Verify that the ports 'f', 'u', 'v', 'w', 'y', and 'h' are closed:

```
# gabconfig -a
GAB Port Memberships
=====
Port a gen 761f03 membership 01
Port b gen 761f08 membership 01
Port d gen 761f02 membership 01
Port o gen 761f01 membership 01
```

- 10 Stop ODM:

```
# /sbin/init.d/odm stop
```

11 Stop LMX:

```
# lmxconfig -U
```

12 Stop VCSMM:

```
# vcsmmconfig -U
```

13 Stop fencing:

```
# /sbin/init.d/vxfen stop
```

14 Stop GAB:

```
# gabconfig -U
```

15 Stop LLT:

```
# lltconfig -U
```

Applying Oracle patches on SF Oracle RAC nodes

Before installing any Oracle RAC patch or patchset software:

- Review the latest information on supported Oracle RAC patches and patchsets:
<http://entsupport.symantec.com/docs/280186>
- You must have installed the base version of the Oracle RAC software.

To apply Oracle patches

- 1 Freeze the service group that contains the Oracle resource and the cssd resource:

```
# hagrps -freeze grp_name
```

- 2 Install the patches or patchsets required for your Oracle RAC installation. See the Oracle documentation that accompanies the patch or patchset.

If the patch is for Oracle database software, after applying the patch, before starting the database or database instance, relink the SF Oracle RAC libraries with Oracle libraries.

- 3 Unfreeze the service group that contains the Oracle resource and the cssd resource:

```
# hagrps -unfreeze grp_name
```

Installing Veritas Volume Manager, Veritas File System, or ODM patches on SF Oracle RAC nodes

Perform the steps on each node in the cluster to install Veritas Volume Manager, Veritas File System, or ODM patches.

To install Veritas Volume Manager, Veritas File System, or ODM patches on SF Oracle RAC nodes

- 1 Log in as the root user.
- 2 Stop SF Oracle RAC on each node.
See [“Stopping SF Oracle RAC using the SF Oracle RAC installer”](#) on page 99.
- 3 Install the VxVM, VxFS, or ODM patch as described in the corresponding patch documentation.
- 4 If there are applications that are not managed by VCS, start the applications manually using native application commands.

Applying operating system updates on SF Oracle RAC nodes

If you need to apply updates to the base version of the operating system, perform the steps in this section on each node of the cluster, one node at a time.

To apply operating system updates

- 1 Log in to the node as the root user and change to `/opt/VRTS/install` directory:

```
# cd /opt/VRTS/install
```

- 2 Take the VCS service groups offline:

```
# hagrps -offline grp_name -sys node_name
```

- 3 Stop SF Oracle RAC:

```
# ./installsfrac<version> -stop
```

Where `<version>` is the specific release version.

- 4 Upgrade the operating system.
See the operating system documentation.

- 5 If the node is not rebooted after the operating system upgrade, reboot the node:

```
# shutdown -r now
```

- 6 Repeat all the steps on each node in the cluster.

Determining the LMX traffic for each database in an SF Oracle RAC cluster

Use the `ltxdbstat` utility to determine the LMX bandwidth used for database traffic for each database. The utility is located at `/sbin/ltxdbstat`.

The utility reports the following information:

- Status of the LMX protocol
- The LMX port and buffer traffic received and transmitted at periodic intervals in packets and kilobytes for each database instance. Specify more than one database when you want to compare database traffic between multiple databases.
- The LMX port and buffer traffic received and transmitted at periodic intervals in packets and kilobytes for a database process.

Note: LMX statistics collection is disabled when you run the command for the first time to view the statistics for a specific database instance or database process. Run the following command to enable LMX statistics collection:

```
# ltxdbstat -z 1
```

Since the utility collects the LMX statistics from `ltxstat` and compiles them for a database, the total processing time may exceed the reported time. Moreover, the percentage utilization reported for a database may vary slightly if some transient Oracle processes consumed LMX bandwidth intermittently along with the instance-specific processes but did not use it at the time of statistics collection.

For more details on the command, see `ltxdbstat (1M)`.

The format of the command is as follows:

```
# ltxdbstat [-v] [-?|-h] [-p <pid>] [-d <dbname>] \  
[interval [count]]  
  
# ltxdbstat -z <0/1>
```

Where:

-v	Displays a verbose output of LMX traffic over the last 1, 10 and 30 second intervals.
-p <i>pid</i>	Displays the statistics for a database process. You need to specify the process ID of the process
-d <i>dbname</i>	Displays the statistics for a database. Specify more than one database when you want to compare database traffic between multiple databases.
-z	0: Stops the collection of LMX statistics. 1: Starts the collection of LMX statistics.
interval	Indicates the period of time in seconds over which LMX statistics is gathered for a database . The default value is 0.
count	Indicates the number of times LMX statistics is gathered for a database. The default value is 1.

Note: Make sure that LMX statistics collection is enabled, otherwise the command fails with the error "LMX stats collection not enabled".

[Table 2-1](#) lists the usage of the command in various scenarios.

Table 2-1 Using lmxdbstat utility to view LMX traffic for databases

Scenario	Command
To view the LMX traffic for all databases	# <code>lmxdbstat</code>
To view the LMX traffic for a particular database instance	# <code>lmxdbstat -d db_name</code>
To compare the traffic between multiple databases	# <code>lmxdbstat -d db_name1 db_name2</code>
To view the LMX traffic for a particular database process	# <code>lmxdbstat -p pid</code>

Table 2-1 Using lmxdbstat utility to view LMX traffic for databases (*continued*)

Scenario	Command
To collect the statistics for a particular interval or frequency for a particular database or all databases	<pre># lmxdbstat interval count</pre> <p>For example, to gather LMX statistics for all databases, 3 times, each for an interval of 10 seconds:</p> <pre># lmxdbstat 10 3</pre>

Adding storage to an SF Oracle RAC cluster

You can add storage to an SF Oracle RAC cluster in the following ways:

Add a disk to a disk group

Use the `vxchg` command to add a disk to a disk group.

See the `vxchg (1M)` manual page for information on various options.

See [“To add storage to an SF Oracle RAC cluster by adding a disk to a disk group”](#) on page 106.

Extend the volume space on a disk group

Use the `vxresize` command to change the length of a volume containing a file system. It automatically locates available disk space on the specified volume and frees up unused space to the disk group for later use.

See the `vxresize (1M)` manual page for information on various options.

See [“To add storage to an SF Oracle RAC cluster by extending the volume space on a disk group”](#) on page 107.

To add storage to an SF Oracle RAC cluster by adding a disk to a disk group

- ◆ Add a disk to the disk group:

```
# vxchg -g dg_name adddisk disk_name
```

To add storage to an SF Oracle RAC cluster by extending the volume space on a disk group

- 1 Determine the length by which you can increase an existing volume.

```
# vxassist [-g diskgroup] maxsize
```

For example, to determine the maximum size the volume `oradata_vol` in the disk group `oradata_dg` can grow, given its attributes and free storage available:

```
# vxassist -g oradata_dg maxsize
```

- 2 Extend the volume, as required. You can extend an existing volume to a certain length by specifying the new size of the volume (the new size must include the additional space you plan to add). You can also extend a volume by a certain length by specifying the additional amount of space you want to add to the existing volume.

To extend a volume to a certain length For example, to extend the volume `oradata_vol` of size 10 GB in the disk group `oradata_dg` to 30 GB:

```
# vxresize -g oradata_dg \  
oradata_vol 30g
```

To extend a volume by a certain length For example, to extend the volume `oradata_vol` of size 10 GB in the disk group `oradata_dg` by 10 GB:

```
# vxresize -g oradata_dg \  
oradata_vol +10g
```

Recovering from storage failure

Veritas Volume Manager (VxVM) protects systems from disk and other hardware failures and helps you to recover from such events. Recovery procedures help you prevent loss of data or system access due to disk and other hardware failures.

For information on various failure and recovery scenarios, see the *Veritas Volume Manager Troubleshooting Guide*.

Backing up and restoring Oracle database using Symantec NetBackup

NetBackup integrates the database backup and recovery capabilities of the Oracle Recovery Manager (RMAN) with the backup and recovery management capabilities

of NetBackup. NetBackup for Oracle also lets you export and import Oracle data in XML format for long-term archival and retrieval.

See the *Symantec NetBackup for Oracle Administrator's Guide*.

Enhancing the performance of SF Oracle RAC clusters

The main components of clustering that impact the performance of an SF Oracle RAC cluster are:

- Kernel components, specifically LLT and GAB
- VCS engine (had)
- VCS agents

Each VCS agent process has two components—the agent framework and the agent functions. The agent framework provides common functionality, such as communication with the HAD, multithreading for multiple resources, scheduling threads, and invoking functions. Agent functions implement functionality that is particular to an agent.

For various options provided by the clustering components to monitor and enhance performance, see the chapter "VCS performance considerations" in the *Veritas Cluster Server Administrator's Guide*.

Veritas Volume Manager can improve system performance by optimizing the layout of data storage on the available hardware.

For more information on tuning Veritas Volume Manager for better performance, see the chapter "Performance monitoring and tuning" in the *Veritas Storage Foundation Administrator's Guide*.

Veritas Volume Replicator Advisor (VRAdvisor) is a planning tool that helps you determine an optimum Veritas Volume Replicator (VVR) configuration.

For installing VRAdvisor and evaluating various parameters using the data collection and data analysis process, see the *Veritas Storage Foundation and High Availability Solutions Replication Administrator's Guide*.

Mounting a snapshot file system for backups increases the load on the system as it involves high resource consumption to perform copy-on-writes and to read data blocks from the snapshot. In such situations, cluster snapshots can be used to do off-host backups. Off-host backups reduce the load of a backup application from the primary server. Overhead from remote snapshots is small when compared to overall snapshot overhead. Therefore, running a backup application by mounting a snapshot from a relatively less loaded node is beneficial to overall cluster performance.

Creating snapshots for offhost processing

You can capture a point-in-time copy of actively changing data at a given instant. You can then perform system backup, upgrade, and other maintenance tasks on the point-in-time copies while providing continuous availability of your critical data. If required, you can offload processing of the point-in-time copies onto another host to avoid contention for system resources on your production server.

See the *Veritas Storage Foundation and High Availability Solutions, Solutions Guide*.

Managing database storage efficiently using SmartTier

SmartTier matches data storage with data usage requirements. After data matching, the data can be relocated based upon data usage and other requirements determined by the storage or database administrator (DBA). The technology enables the database administrator to manage data so that less frequently used data can be moved to slower, less expensive disks. This also permits the frequently accessed data to be stored on faster disks for quicker retrieval.

See the *Veritas Storage Foundation: Storage and Availability Management for Oracle Databases* guide.

Optimizing database storage using Thin Provisioning and SmartMove

Thin Provisioning is a storage array feature that optimizes storage use by automating storage provisioning. With Storage Foundation's SmartMove feature, Veritas File System (VxFS) lets Veritas Volume Manager (VxVM) know which blocks have data. VxVM, which is the copy engine for migration, copies only the used blocks and avoids the empty spaces, thus optimizing thin storage utilization.

See the *Veritas Storage Foundation and High Availability Solutions, Solutions Guide*.

Scheduling periodic health checks for your SF Oracle RAC cluster

You can manually invoke the health check utility at any time or schedule the utility to run as a cron job.

If the health check completes without any warnings, the following message is displayed:

```
Success: The SF Oracle RAC components are working fine on this node.
```

If you manually invoke the health check, run the utility on each node in the cluster. The results of the check are displayed on the console.

To run health checks on a cluster

- 1 Log in as the root user on each node in the cluster.
- 2 Using `vi` or any text editor, modify the health check parameters as required for your installation setup.

Note: You must set the `ORACLE_HOME` and `CRS_HOME` parameters in the configuration file as appropriate for your setup.

```
# cd /opt/VRTSvcs/rac/healthcheck/
# vi healthcheck.cf
```

- 3 Run the health check utility:

```
# ./healthcheck
```

If you want to schedule periodic health checks for your cluster, create a cron job that runs on each node in the cluster. Redirect the health check report to a file.

Verifying the nodes in an SF Oracle RAC cluster

[Table 2-2](#) lists the various options that you can use to periodically verify the nodes in your cluster.

Table 2-2 Options for verifying the nodes in a cluster

Type of check	Description
Symantec Operations Readiness Tools (SORT)	<p>Use the Symantec Operations Readiness Tools (SORT) to evaluate your systems before and after any installation, configuration, upgrade, patch updates, or other routine administrative activities. The utility performs a number of compatibility and operational checks on the cluster that enable you to diagnose and troubleshoot issues in the cluster. The utility is periodically updated with new features and enhancements.</p> <p>For more information and to download the utility, visit http://sort.symantec.com.</p>

Table 2-2 Options for verifying the nodes in a cluster (*continued*)

Type of check	Description
SF Oracle RAC health checks	<p>SF Oracle RAC provides a health check utility that examines the functional health of the components in an SF Oracle RAC cluster. The utility when invoked gathers real-time operational information on cluster components and displays the report on your system console. You must run the utility on each node in the cluster.</p> <p>To schedule periodic health checks:</p> <p>See “Periodic health evaluation of SF Oracle RAC clusters” on page 87.</p>

Administering VCS

This section provides instructions for the following VCS administration tasks:

- Viewing available Veritas devices and drivers
See [“Viewing available Veritas device drivers”](#) on page 112.
- Loading Veritas drivers into memory
See [“Loading Veritas drivers into memory”](#) on page 112.
- Starting and stopping VCS
See [“Starting and stopping VCS”](#) on page 112.
- Environment variables to start and stop VCS modules
See [“Environment variables to start and stop VCS modules”](#) on page 113.
- Adding and removing LLT links
See [“Adding and removing LLT links”](#) on page 115.
- Displaying the cluster details and LLT version for LLT links
See [“Displaying the cluster details and LLT version for LLT links”](#) on page 121.
- Configuring aggregated interfaces under LLT
See [“Configuring aggregated interfaces under LLT”](#) on page 119.
- Configuring destination-based load balancing for LLT
See [“Configuring destination-based load balancing for LLT”](#) on page 122.
- Enabling and disabling intelligent resource monitoring
See [“Enabling and disabling intelligent resource monitoring for agents manually”](#) on page 122.
- Administering the AMF kernel driver
See [“Administering the AMF kernel driver”](#) on page 124.

If you encounter issues while administering VCS, see the troubleshooting section in the *Veritas Cluster Server Administrator's Guide* for assistance.

Viewing available Veritas device drivers

To view the devices that are loaded in memory, run the `kcmodule` command as shown in the following examples.

For example:

If you want to view whether or not the driver 'gab' is loaded in memory:

```
# kcmodule |grep gab
gab                loaded  explicit  auto-loadable, unloadable
```

If you want to view whether or not the 'vx' drivers are loaded in memory:

```
# kcmodule |grep vx

vxdump            static  best
vxfen             loaded  explicit auto-loadable, unloadable
vxfs              unused
vxfs60            static  best      loadable, unloadable
vxglm             loaded  explicit auto-loadable, unloadable
vxgms            loaded  explicit auto-loadable, unloadable
vxportal          unused  auto-loadable, unloadable
vxportal60        static  best      loadable, unloadable
```

Loading Veritas drivers into memory

Under normal operational conditions, you do not need to load Veritas drivers into memory. You might need to load a Veritas driver only if there is a malfunction.

To load the ODM driver into memory, for example:

```
# kcmodule odm=loaded
```

Starting and stopping VCS

To start VCS on each node:

```
# hastart
```

To stop VCS on each node:

```
# hastop -local
```


You can also use the command `hastop -all`; however, make sure that you wait for port 'h' to close before restarting VCS.

Environment variables to start and stop VCS modules

The start and stop environment variables for AMF, LLT, GAB, VxFEN, VCSMM, LMX, and VCS engine define the default VCS behavior to start these modules during system restart or stop these modules during system shutdown.

Note: The startup and shutdown of AMF, LLT, GAB, VxFEN, VCSMM, LMX, and VCS engine are inter-dependent. For a clean startup or shutdown of SF Oracle RAC, you must either enable or disable the startup and shutdown modes for all these modules.

[Table 2-3](#) describes the start and stop variables for VCS.

Table 2-3 Start and stop environment variables for VCS

Environment variable	Definition and default value
AMF_START	Startup mode for the AMF driver. By default, the AMF driver is enabled to start up after a system reboot. This environment variable is defined in the following file: <code>/etc/rc.config.d/amf</code> Default: 1
AMF_STOP	Shutdown mode for the AMF driver. By default, the AMF driver is enabled to stop during a system shutdown. This environment variable is defined in the following file: <code>/etc/rc.config.d/amf</code> Default: 1
LLT_START	Startup mode for LLT. By default, LLT is enabled to start up after a system reboot. This environment variable is defined in the following file: <code>/etc/rc.config.d/lltconf</code> Default: 1

Table 2-3 Start and stop environment variables for VCS (*continued*)

Environment variable	Definition and default value
LLT_STOP	<p>Shutdown mode for LLT. By default, LLT is enabled to stop during a system shutdown.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/lltconf</code></p> <p>Default: 1</p>
GAB_START	<p>Startup mode for GAB. By default, GAB is enabled to start up after a system reboot.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/gabconf</code></p> <p>Default: 1</p>
GAB_STOP	<p>Shutdown mode for GAB. By default, GAB is enabled to stop during a system shutdown.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/gabconf</code></p> <p>Default: 1</p>
VXFEN_START	<p>Startup mode for VxFEN. By default, VxFEN is enabled to start up after a system reboot.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/vxfenconf</code></p> <p>Default: 1</p>
VXFEN_STOP	<p>Shutdown mode for VxFEN. By default, VxFEN is enabled to stop during a system shutdown.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/vxfenconf</code></p> <p>Default: 1</p>
VCSMM_START	<p>Startup mode for VCSMM. By default, VCSMM is enabled to start up after a system reboot.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/vcsmmconf</code></p> <p>Default: 1</p>

Table 2-3 Start and stop environment variables for VCS (*continued*)

Environment variable	Definition and default value
VCSMM_STOP	<p>Shutdown mode for VCSMM. By default, VCSMM is enabled to stop during a system shutdown.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/vcsmmconf</code></p> <p>Default: 1</p>
LMX_START	<p>Startup mode for LMX. By default, LMX is enabled to start up after a system reboot.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/lmxconf</code></p> <p>Default: 1</p>
LMX_STOP	<p>Shutdown mode for LMX. By default, LMX is enabled to stop during a system shutdown.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/lmxconf</code></p> <p>Default: 1</p>
VCS_START	<p>Startup mode for VCS engine. By default, VCS engine is enabled to start up after a system reboot.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/vcsconf</code></p> <p>Default: 1</p>
VCS_STOP	<p>Shutdown mode for VCS engine. By default, VCS engine is enabled to stop during a system shutdown.</p> <p>This environment variable is defined in the following file:</p> <p><code>/etc/rc.config.d/vcsconf</code></p> <p>Default: 1</p>

Adding and removing LLT links

In an SF Oracle RAC cluster, Oracle Clusterware heartbeat link must be configured as an LLT link. If Oracle Clusterware and LLT use different links for their communication, then the membership change between VCS and Oracle Clusterware

is not coordinated correctly. For example, if only the Oracle Clusterware links are down, Oracle Clusterware kills one set of nodes after the expiry of the `css-misscount` interval and initiates the Oracle Clusterware and database recovery, even before CVM and CFS detect the node failures. This uncoordinated recovery may cause data corruption.

If you need additional capacity for Oracle communication on your private interconnects, you can add LLT links. The network IDs of the interfaces connected to the same physical network must match. The interfaces specified in the PrivNIC or MultiPrivNIC configuration must be exactly the same in name and total number as those which have been used for LLT configuration.

You can use the `lltconfig` command to add or remove LLT links when LLT is running.

LLT links can be added or removed while clients are connected.

See the `lltconfig(1M)` manual page for more details.

Note: When you add or remove LLT links, you need not shut down GAB or the high availability daemon, `had`. Your changes take effect immediately, but are lost on the next restart. For changes to persist, you must also update the `/etc/llttab` file.

To add LLT links

- 1 Depending on the LLT link type, run the following command to add an LLT link:

- For ether link type:

```
# lltconfig -t devtag -d device  
[-b ether ] [-s SAP] [-m mtu]
```

- For UDP link type:

```
# lltconfig -t devtag -d device  
-b udp [-s port] [-m mtu]  
-I IPaddr -B bcst
```

- For UDP6 link type:

```
# lltconfig -t devtag -d device  
-b udp6 [-s port] [-m mtu]  
-I IPaddr [-B mcast]
```

Where:

devtag	Tag to identify the link
device	Network device path of the interface For link type ether, the path is followed by a colon (:) and an integer which specifies the unit or PPA used by LLT to attach. For link types udp and udp6, the device is the udp and udp6 device path respectively.
bcast	Broadcast address for the link type udp
mcast	Multicast address for the link type udp6
IPaddr	IP address for link types udp and udp6
SAP	SAP to bind on the network links for link type ether
port	Port for link types udp and udp6
mtu	Maximum transmission unit to send packets on network links

For example:

■ For ether link type:

```
# lltnconfig -t lan3 -d /dev/lan:3 -s 0xcafe -m 1500
```

■ For UDP link type:

```
# lltnconfig -t link1 -d /dev/udp -b udp  
-I 192.1.2.255 -B 192.1.2.255
```

■ For UDP6 link type:

```
# lltnconfig -t link1 -d /dev/udp6  
-b udp6 -I 2000::1
```

Note: If you want the addition of LLT links to be persistent after reboot, then you must edit the `/etc/lltab` with LLT entries.

- 2 If you want to configure the link under PrivNIC or MultiPrivNIC as a failover target in the case of link failure, modify the PrivNIC or MultiPrivNIC configuration as follows:

```
# haconf -makerw
# hares -modify resource_name Device device
    device_id [-sys hostname]
# haconf -dump -makero
```

The following is an example of configuring the link under PrivNIC.

Assuming that you have two LLT links configured under PrivNIC as follows:

```
PrivNIC ora_priv (
    Critical = 0
    Device@sys1 = { lan1 = 0, lan2 = 1 }
    Device@sys2 = { lan1 = 0, lan2 = 1 }
    Address@sys1 = "192.168.12.1"
    Address@sys2 = "192.168.12.2"
    NetMask = "255.255.255.0"
)
```

To configure the new LLT link under PrivNIC, run the following commands:

```
# haconf -makerw
# hares -modify ora_priv Device lan1 0 lan2 1 lan3 2 -sys sys1
# hares -modify ora_priv Device lan1 0 lan2 1 lan3 2 -sys sys2
# haconf -dump -makero
```

The updated PrivNIC resource now resembles:

```
PrivNIC ora_priv (
    Critical = 0
    Device@sys1 = { lan1 = 0, lan2 = 1, lan3 = 2 }
    Device@sys2 = { lan1 = 0, lan2 = 1, lan3 = 2 }
    Address@sys1 = "192.168.12.1"
    Address@sys2 = "192.168.12.2"
    NetMask = "255.255.255.0"
)
```

To remove an LLT link

- 1 If the link you want to remove is configured as a PrivNIC or MultiPrivNIC resource, you need to modify the resource configuration before removing the link.

If you have configured the link under PrivNIC or MultiPrivNIC as a failover target in the case of link failure, modify the PrivNIC or MultiPrivNIC configuration as follows:

```
# haconf -makerw
# hares -modify resource_name Device link_name \
device_id [-sys hostname]
# haconf -dump -makero
```

For example, if the links lan1, lan2, and lan3 were configured as PrivNIC resources, and you want to remove lan3:

```
# haconf -makerw
# hares -modify ora_priv Device lan1
    0 \
lan2
    1
```

where lan1 and lan2 are the links that you want to retain in your cluster.

```
# haconf -dump -makero
```

- 2 Run the following command to remove a network link that is configured under LLT:

```
# lltconfig -u devtag
```

Configuring aggregated interfaces under LLT

If you want to configure LLT to use aggregated interfaces after installing and configuring VCS, you can use one of the following approaches:

- Edit the `/etc/llttab` file
This approach requires you to stop LLT. The aggregated interface configuration is persistent across reboots.
- Run the `lltconfig` command
This approach lets you configure aggregated interfaces on the fly. However, the changes are not persistent across reboots.

To configure aggregated interfaces under LLT by editing the /etc/llttab file

- 1 If LLT is running, stop LLT after you stop the other dependent modules.

```
# /sbin/init.d/llt stop
```

See “Starting or stopping SF Oracle RAC on each node” on page 97.

- 2 Add the following entry to the /etc/llttab file to configure an aggregated interface.

```
link tag device_name systemid_range link_type sap mtu_size
```

tag	Tag to identify the link
device_name	Network device path of the aggregated interface The path is followed by a colon (:) and an integer which specifies the unit or PPA used by LLT to attach.
systemid_range	Range of systems for which the command is valid. If the link command is valid for all systems, specify a dash (-).
link_type	The link type must be ether.
sap	SAP to bind on the network links. Default is 0xcaff.
mtu_size	Maximum transmission unit to send packets on network links

- 3 Restart LLT for the changes to take effect. Restart the other dependent modules that you stopped in step 1.

```
# /sbin/init.d/llt start
```

See “Starting or stopping SF Oracle RAC on each node” on page 97.

To configure aggregated interfaces under LLT using the `lltconfig` command

- ◆ When LLT is running, use the following command to configure an aggregated interface:

```
lltconfig -t devtag -d device  
[-b linktype ] [-s SAP] [-m mtu]
```

<code>devtag</code>	Tag to identify the link
<code>device</code>	Network device path of the aggregated interface The path is followed by a colon (:) and an integer which specifies the unit or PPA used by LLT to attach.
<code>link_type</code>	The link type must be ether.
<code>sap</code>	SAP to bind on the network links. Default is 0xcale.
<code>mtu_size</code>	Maximum transmission unit to send packets on network links

See the `lltconfig(1M)` manual page for more details.

You need not reboot after you make this change. However, to make these changes persistent across reboot, you must update the `/etc/llttab` file.

See [“To configure aggregated interfaces under LLT by editing the `/etc/llttab` file”](#) on page 120.

Displaying the cluster details and LLT version for LLT links

You can use the `lltdump` command to display the LLT version for a specific LLT link. You can also display the cluster ID and node ID details.

See the `lltdump(1M)` manual page for more details.

To display the cluster details and LLT version for LLT links

- ◆ Run the following command to display the details:

```
# /opt/VRTSllt/lltdump -D -f link
```

For example, if lan3 is connected to galaxy, then the command displays a list of all cluster IDs and node IDs present on the network link lan3.

```
# /opt/VRTSllt/lltdump -D -f /dev/lan:3
```

```
lltdump : Configuration:

device : lan3

sap : 0xcafe
promisc sap : 0
promisc mac : 0
cidsnoop : 1
=== Listening for LLT packets ===
cid nid vmaj vmin
3456 1 5 0
3456 3 5 0
83 0 4 0
27 1 3 7
3456 2 5 0
```

Configuring destination-based load balancing for LLT

Destination-based load balancing for LLT is turned off by default. Symantec recommends destination-based load balancing when the cluster setup has more than two nodes and more active LLT ports.

See [“About Low Latency Transport \(LLT\)”](#) on page 25.

To configure destination-based load balancing for LLT

- ◆ Run the following command to configure destination-based load balancing:

```
lltconfig -F linkburst:0
```

Enabling and disabling intelligent resource monitoring for agents manually

Review the following procedures to enable or disable intelligent resource monitoring manually. The intelligent resource monitoring feature is enabled by

default. The IMF resource type attribute determines whether an IMF-aware agent must perform intelligent resource monitoring.

See “[About resource monitoring](#)” on page 36.

To enable intelligent resource monitoring

- 1 Make the VCS configuration writable.

```
# haconf -makerw
```

- 2 Run the following command to enable intelligent resource monitoring.

- To enable intelligent monitoring of offline resources:

```
# hatype -modify resource_type IMF -update Mode 1
```

- To enable intelligent monitoring of online resources:

```
# hatype -modify resource_type IMF -update Mode 2
```

- To enable intelligent monitoring of both online and offline resources:

```
# hatype -modify resource_type IMF -update Mode 3
```

- 3 If required, change the values of the MonitorFreq key and the RegisterRetryLimit key of the IMF attribute.

See the *Veritas Cluster Server Bundled Agents Reference Guide* for agent-specific recommendations to set these attributes.

Review the agent-specific recommendations in the attribute definition tables to set these attribute key values.

- 4 Save the VCS configuration.

```
# haconf -dump -makero
```

- 5 Restart the agent. Run the following commands on each node.

```
# haagent -stop agent_name -force -sys sys_name
```

```
# haagent -start agent_name -sys sys_name
```

To disable intelligent resource monitoring

- 1 Make the VCS configuration writable.

```
# haconf -makerw
```

- 2 To disable intelligent resource monitoring for all the resources of a certain type, run the following command:

```
# hatype -modify resource_type IMF -update Mode 0
```

- 3 To disable intelligent resource monitoring for a specific resource, run the following command:

```
# hares -override resource_name IMF  
# hares -modify resource_name IMF -update Mode 0
```

- 4 Save the VCS configuration.

```
# haconf -dump -makero
```

Note: VCS provides haimfconfig script to enable or disable the IMF functionality for agents. You can use the script with VCS in running or stopped state. Use the script to enable or disable IMF for the IMF-aware bundled agents, enterprise agents, and custom agents.

Administering the AMF kernel driver

Review the following procedures to start, stop, or unload the AMF kernel driver.

See [“About the IMF notification module”](#) on page 35.

See [“Environment variables to start and stop VCS modules”](#) on page 113.

To start the AMF kernel driver

- 1 Set the value of the AMF_START variable to 1 in the following file, if the value is not already 1:

```
# /etc/rc.config.d/amf
```

- 2 Start the AMF kernel driver. Run the following command:

```
# /sbin/init.d/amf start
```

To stop the AMF kernel driver

- 1 Set the value of the AMF_START variable to 0 in the following file, if the value is not already 0:

```
# /etc/rc.config.d/amf
```

- 2 Stop the AMF kernel driver. Run the following command:

```
# /sbin/init.d/amf stop
```

To unload the AMF kernel driver

- 1 If agent downtime is not a concern, use the following steps to unload the AMF kernel driver:

- Stop the agents that are registered with the AMF kernel driver.
The `amfstat` command output lists the agents that are registered with AMF under the Registered Reapers section.
See the `amfstat` manual page.
- Stop the AMF kernel driver.
See [“To stop the AMF kernel driver”](#) on page 125.
- Start the agents.

- 2 If you want minimum downtime of the agents, use the following steps to unload the AMF kernel driver:

- Run the following command to disable the AMF driver even if agents are still registered with it.

```
# amfconfig -Uof
```

- Stop the AMF kernel driver.
See [“To stop the AMF kernel driver”](#) on page 125.

Administering I/O fencing

This section describes I/O fencing and provides instructions for common I/O fencing administration tasks.

- About administering I/O fencing
See [“About administering I/O fencing”](#) on page 126.
- About `vxfcntlsthdw` utility
See [“About the `vxfcntlsthdw` utility”](#) on page 127.

- About vxfenadm utility
See [“About the vxfenadm utility”](#) on page 134.
- About vxfenclearpre utility
See [“About the vxfenclearpre utility”](#) on page 139.
- About vxfenswap utility
See [“About the vxfenswap utility”](#) on page 142.

If you encounter issues while administering I/O fencing, refer to the troubleshooting section for assistance.

See [“Troubleshooting I/O fencing”](#) on page 199.

See [“About administering I/O fencing”](#) on page 126.

About administering I/O fencing

The I/O fencing feature provides the following utilities that are available through the `VRTSvxfen` depot:

<code>vxfentsthdw</code>	Tests SCSI-3 functionality of the disks for I/O fencing See “About the vxfentsthdw utility” on page 127.
<code>vxfenconfig</code>	Configures and unconfigures I/O fencing Lists the coordination points used by the vxfen driver.
<code>vxfenadm</code>	Displays information on I/O fencing operations and manages SCSI-3 disk registrations and reservations for I/O fencing See “About the vxfenadm utility” on page 134.
<code>vxfenclearpre</code>	Removes SCSI-3 registrations and reservations from disks See “About the vxfenclearpre utility” on page 139.
<code>vxfenswap</code>	Replaces coordination points without stopping I/O fencing See “About the vxfenswap utility” on page 142.
<code>vx fendisk</code>	Generates the list of paths of disks in the diskgroup. This utility requires that Veritas Volume Manager is installed and configured.

The I/O fencing commands reside in the `/opt/VRTS/bin|grep -i vxfen` folder. Make sure you added this folder path to the PATH environment variable.

Refer to the corresponding manual page for more information on the commands.

About the vxfcntlsthdw utility

You can use the `vxfcntlsthdw` utility to verify that shared storage arrays to be used for data support SCSI-3 persistent reservations and I/O fencing. During the I/O fencing configuration, the testing utility is used to test a single disk. The utility has other options that may be more suitable for testing storage devices in other configurations. You also need to test coordinator disk groups.

See *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* to set up I/O fencing.

The utility, which you can run from one system in the cluster, tests the storage used for data by setting and verifying SCSI-3 registrations on the disk or disks you specify, setting and verifying persistent reservations on the disks, writing data to the disks and reading it, and removing the registrations from the disks.

Refer also to the `vxfcntlsthdw(1M)` manual page.

About general guidelines for using the vxfcntlsthdw utility

Review the following guidelines to use the `vxfcntlsthdw` utility:

- The utility requires two systems connected to the shared storage.

Caution: The tests overwrite and destroy data on the disks, unless you use the `-r` option.

- The two nodes must have `ssh` (default) or `remsh` communication. If you use `remsh`, launch the `vxfcntlsthdw` utility with the `-n` option.
After completing the testing process, you can remove permissions for communication and restore public network connections.
- To ensure both systems are connected to the same disk during the testing, you can use the `vxflenadm -i diskpath` command to verify a disk's serial number. See [“Verifying that the nodes see the same disk”](#) on page 139.
- For disk arrays with many disks, use the `-m` option to sample a few disks before creating a disk group and using the `-g` option to test them all.
- The utility indicates a disk can be used for I/O fencing with a message resembling:

```
The disk /dev/vx/rmp/ctl1d0 is ready to be configured for
I/O Fencing on node sys1
```

If the utility does not show a message stating a disk is ready, verification has failed.

- If the disk you intend to test has existing SCSI-3 registration keys, the test issues a warning before proceeding.

About the vxfcntlshdw command options

Table 2-4 describes the methods that the utility provides to test storage devices.

Table 2-4 vxfcntlshdw options

vxfcntlshdw option	Description	When to use
-n	Utility uses remsh for communication.	Use when remsh is used for communication.
-r	Non-destructive testing. Testing of the disks for SCSI-3 persistent reservations occurs in a non-destructive way; that is, there is only testing for reads, not writes. May be used with -m, -f, or -g options.	Use during non-destructive testing. See “Performing non-destructive testing on the disks using the -r option” on page 131.
-t	Testing of the return value of SCSI TEST UNIT (TUR) command under SCSI-3 reservations. A warning is printed on failure of TUR testing.	When you want to perform TUR testing.
-d	Use DMP devices. May be used with -c or -g options.	By default, the script picks the DMP paths for disks in the diskgroup. If you want the script to use the raw paths for disks in the diskgroup, use the -w option.
-w	Use raw devices. May be used with -c or -g options.	With the -w option, the script picks the raw paths for disks in the diskgroup. By default, the script uses the -d option to pick up the DMP paths for disks in the disk group.
-c	Utility tests the coordinator disk group prompting for systems and devices, and reporting success or failure.	For testing disks in coordinator disk group. See “Testing the coordinator disk group using vxfcntlshdw -c option” on page 129.

Table 2-4 vxfststhdw options (*continued*)

vxfststhdw option	Description	When to use
-m	Utility runs manually, in interactive mode, prompting for systems and devices, and reporting success or failure. May be used with -r and -t options. -m is the default option.	For testing a few disks or for sampling disks in larger arrays. See “Testing the shared disks using the vxfststhdw -m option” on page 131.
-f <i>filename</i>	Utility tests system/device combinations listed in a text file. May be used with -r and -t options.	For testing several disks. See “Testing the shared disks listed in a file using the vxfststhdw -f option” on page 133.
-g <i>disk_group</i>	Utility tests all disk devices in a specified disk group. May be used with -r and -t options.	For testing many disks and arrays of disks. Disk groups may be temporarily created for testing purposes and destroyed (ungrouped) after testing. See “Testing all the disks in a disk group using the vxfststhdw -g option” on page 133.

Testing the coordinator disk group using vxfststhdw -c option

Use the vxfststhdw utility to verify disks are configured to support I/O fencing. In this procedure, the vxfststhdw utility tests the three disks one disk at a time from each node.

The procedure in this section uses the following disks for example:

- From the node sys1, the disks are seen as /dev/vx/rdmp/c1t1d0, /dev/vx/rdmp/c2t1d0, and /dev/vx/rdmp/c3t1d0.
- From the node sys2, the same disks are seen as /dev/vx/rdmp/c4t1d0, /dev/vx/rdmp/c5t1d0, and /dev/vx/rdmp/c6t1d0.

Note: To test the coordinator disk group using the vxfststhdw utility, the utility requires that the coordinator disk group, vxencoorddg, be accessible from two nodes.

To test the coordinator disk group using `vxfcntlshdw -c`

- 1 Use the `vxfcntlshdw` command with the `-c` option. For example:

```
# vxfcntlshdw -c vxfencoorddg
```

- 2 Enter the nodes you are using to test the coordinator disks:

```
Enter the first node of the cluster: sys1
```

```
Enter the second node of the cluster: sys2
```

- 3 Review the output of the testing process for both nodes for all disks in the coordinator disk group. Each disk should display output that resembles:

```
ALL tests on the disk /dev/vx/rmp/c1t1d0 have PASSED.  
The disk is now ready to be configured for I/O Fencing on node  
sys1 as a COORDINATOR DISK.
```

```
ALL tests on the disk /dev/vx/rmp/c4t1d0 have PASSED.  
The disk is now ready to be configured for I/O Fencing on node  
sys2 as a COORDINATOR DISK.
```

- 4 After you test all disks in the disk group, the `vxfencoorddg` disk group is ready for use.

Removing and replacing a failed disk

If a disk in the coordinator disk group fails verification, remove the failed disk or LUN from the `vxfencoorddg` disk group, replace it with another, and retest the disk group.

To remove and replace a failed disk

- 1 Use the `vxdiskadm` utility to remove the failed disk from the disk group.
Refer to the *Veritas Storage Foundation Administrator's Guide*.
- 2 Add a new disk to the node, initialize it, and add it to the coordinator disk group.
See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for instructions to initialize disks for I/O fencing and to set up coordinator disk groups.
If necessary, start the disk group.
See the *Veritas Storage Foundation Administrator's Guide* for instructions to start the disk group.
- 3 Retest the disk group.
See [“Testing the coordinator disk group using `vxfcntlsthdw -c` option”](#) on page 129.

Performing non-destructive testing on the disks using the `-r` option

You can perform non-destructive testing on the disk devices when you want to preserve the data.

To perform non-destructive testing on disks

- ◆ To test disk devices containing data you want to preserve, you can use the `-r` option with the `-m`, `-f`, or `-g` options.

For example, to use the `-m` option and the `-r` option, you can run the utility as follows:

```
# vxfcntlsthdw -rm
```

When invoked with the `-r` option, the utility does not use tests that write to the disks. Therefore, it does not test the disks for all of the usual conditions of use.

Testing the shared disks using the `vxfcntlsthdw -m` option

Review the procedure to test the shared disks. By default, the utility uses the `-m` option.

This procedure uses the `/dev/vx/rdmp/c1t1d0` disk in the steps.

If the utility does not show a message stating a disk is ready, verification has failed. Failure of verification can be the result of an improperly configured disk array. It can also be caused by a bad disk.

If the failure is due to a bad disk, remove and replace it. The `vxfcntlshdw` utility indicates a disk can be used for I/O fencing with a message resembling:

```
The disk /dev/vx/rdmp/c1t1d0 is ready to be configured for  
I/O Fencing on node sys1
```

Note: For A/P arrays, run the `vxfcntlshdw` command only on active enabled paths.

To test disks using `vxfcntlshdw` script

- 1 Make sure system-to-system communication is functioning properly.
- 2 From one node, start the utility.

```
# vxfcntlshdw [-n]
```

- 3 After reviewing the overview and warning that the tests overwrite data on the disks, confirm to continue the process and enter the node names.

```
***** WARNING!!!!!!!!!! *****  
THIS UTILITY WILL DESTROY THE DATA ON THE DISK!!  
  
Do you still want to continue : [y/n] (default: n) y  
Enter the first node of the cluster: sys1  
Enter the second node of the cluster: sys2
```

- 4 Enter the names of the disks you are checking. For each node, the disk may be known by the same name:

```
Enter the disk name to be checked for SCSI-3 PGR on node  
sys1 in the format: /dev/vx/rdmp/cxtxdx  
/dev/vx/rdmp/c2t13d0  
Enter the disk name to be checked for SCSI-3 PGR on node  
sys2 in the format: /dev/vx/rdmp/cxtxdx  
Make sure it's the same disk as seen by nodes sys1 and sys2  
/dev/vx/rdmp/c2t13d0
```

If the serial numbers of the disks are not identical, then the test terminates.

- 5 Review the output as the utility performs the checks and report its activities.

- 6 If a disk is ready for I/O fencing on each node, the utility reports success:

```
ALL tests on the disk /dev/vx/rdmp/c1t1d0 have PASSED
The disk is now ready to be configured for I/O Fencing on node
sys1
...
Removing test keys and temporary files, if any ...
.
.
```

- 7 Run the `vxfcntlsthaw` utility for each disk you intend to verify.

Testing the shared disks listed in a file using the `vxfcntlsthaw -f` option

Use the `-f` option to test disks that are listed in a text file. Review the following example procedure.

To test the shared disks listed in a file

- 1 Create a text file `disks_test` to test two disks shared by systems `sys1` and `sys2` that might resemble:

```
sys1 /dev/vx/rdmp/c2t2d1 sys2 /dev/vx/rdmp/c3t2d1
sys1 /dev/vx/rdmp/c2t2d1 sys2 /dev/vx/rdmp/c3t2d1
```

where the first disk is listed in the first line and is seen by `sys1` as `/dev/vx/rdmp/c2t2d1` and by `sys2` as `/dev/vx/rdmp/c3t2d1`. The other disk, in the second line, is seen as `/dev/vx/rdmp/c2t2d2` from `sys1` and `/dev/vx/rdmp/c3t2d2` from `sys2`. Typically, the list of disks could be extensive.

- 2 To test the disks, enter the following command:

```
# vxfcntlsthaw -f disks_test
```

The utility reports the test results one disk at a time, just as for the `-m` option.

Testing all the disks in a disk group using the `vxfcntlsthaw -g` option

Use the `-g` option to test all disks within a disk group. For example, you create a temporary disk group consisting of all disks in a disk array and test the group.

Note: Do not import the test disk group as shared; that is, do not use the `-s` option with the `vxdbg import` command.

After testing, destroy the disk group and put the disks into disk groups as you need.

To test all the disks in a diskgroup

- 1 Create a diskgroup for the disks that you want to test.
- 2 Enter the following command to test the diskgroup test_disks_dg:

```
# vxfentsthdw -g test_disks_dg
```

The utility reports the test results one disk at a time.

Testing a disk with existing keys

If the utility detects that a coordinator disk has existing keys, you see a message that resembles:

```
There are Veritas I/O fencing keys on the disk. Please make sure that I/O fencing is shut down on all nodes of the cluster before continuing.
```

```
***** WARNING!!!!!!!!!! *****
```

```
THIS SCRIPT CAN ONLY BE USED IF THERE ARE NO OTHER ACTIVE NODES IN THE CLUSTER! VERIFY ALL OTHER NODES ARE POWERED OFF OR INCAPABLE OF ACCESSING SHARED STORAGE.
```

```
If this is not the case, data corruption will result.
```

```
Do you still want to continue : [y/n] (default: n) y
```

The utility prompts you with a warning before proceeding. You may continue as long as I/O fencing is not yet configured.

About the vxfenadm utility

Administrators can use the vxfenadm command to troubleshoot and test fencing configurations.

The command's options for use by administrators are as follows:

-s read the keys on a disk and display the keys in numeric, character, and node format

Note: The -g and -G options are deprecated. Use the -s option.

-i read SCSI inquiry information from device

-m	register with disks
-n	make a reservation with disks
-p	remove registrations made by other systems
-r	read reservations
-x	remove registrations

Refer to the `vxfsenadm(1m)` manual page for a complete list of the command options.

About the I/O fencing registration key format

The keys that the `vxfsen` driver registers on the data disks and the coordinator disks consist of eight bytes. The key format is different for the coordinator disks and data disks.

The key format of the coordinator disks is as follows:

Byte	0	1	2	3	4	5	6	7
Value	V	F	cID 0x	cID 0x	cID 0x	cID 0x	nID 0x	nID 0x

where:

- VF is the unique identifier that carves out a namespace for the keys (consumes two bytes)
- cID 0x is the LLT cluster ID in hexadecimal (consumes four bytes)
- nID 0x is the LLT node ID in hexadecimal (consumes two bytes)

The `vxfsen` driver uses this key format in both sybase mode of I/O fencing.

The key format of the data disks that are configured as failover disk groups under VCS is as follows:

Byte	0	1	2	3	4	5	6	7
Value	A+nID	V	C	S				

where nID is the LLT node ID

For example: If the node ID is 1, then the first byte has the value as B ('A' + 1 = B).

The key format of the data disks configured as parallel disk groups under CVM is as follows:

Byte	0	1	2	3	4	5	6	7
Value	A+nID	P	G	R	DGcount	DGcount	DGcount	DGcount

where DGcount is the count of disk group in the configuration (consumes four bytes).

By default, CVM uses unique fencing key for each disk group. However, some arrays have a restriction on the total number of unique keys that can be registered. In such cases, you can use the `same_key_for_alldgs` tunable parameter to change the default behavior. The default value of the parameter is off. If your configuration hits the storage array limit on total number of unique keys, you can turn the value on using the `vxdefault` command as follows:

```
# vxdefault set same_key_for_alldgs on
# vxdefault list
KEYWORD                CURRENT-VALUE  DEFAULT-VALUE
...
same_key_for_alldgs    on              off
...
```

If the tunable is changed to 'on', all subsequent keys that the CVM generates on disk group imports or creates have '0000' as their last four bytes (DGcount is 0). You must deport and re-import all the disk groups that are already imported for the changed value of the `same_key_for_alldgs` tunable to take effect.

Displaying the I/O fencing registration keys

You can display the keys that are currently assigned to the disks using the `vxfenadm` command.

The variables such as `disk_7`, `disk_8`, and `disk_9` in the following procedure represent the disk names in your setup.

To display the I/O fencing registration keys

- 1 To display the key for the disks, run the following command:

```
# vxfenadm -s disk_name
```

For example:

- To display the key for the coordinator disk `/dev/vx/rdmp/c1t1d0` from the system with node ID 1, enter the following command:

```
# vxfenadm -s /dev/vx/rdmp/c1t1d0
key[1]:
[Numeric Format]: 86,70,68,69,69,68,48,48
```



```
[Character Format]: VFDEED00
* [Node Format]: Cluster ID: 57069 Node ID: 0 Node Name: sys1
```

The `-s` option of `vxfenadm` displays all eight bytes of a key value in three formats. In the numeric format,

- The first two bytes, represent the identifier VF, contains the ASCII value 86, 70.
 - The next four bytes contain the ASCII value of the cluster ID 57069 encoded in hex (0xDEED) which are 68, 69, 69, 68.
 - The remaining bytes contain the ASCII value of the node ID 0 (0x00) which are 48, 48. Node ID 1 would be 01 and node ID 10 would be 0A.
- An asterisk before the Node Format indicates that the `vxfenadm` command is run from the node of a cluster where LLT is configured and running.

- To display the keys on a CVM parallel disk group:

```
# vxfenadm -s /dev/vx/rdmp/disk_7

Reading SCSI Registration Keys...

Device Name: /dev/vx/rdmp/disk_7
Total Number Of Keys: 1
key[0]:
  [Numeric Format]: 66,80,71,82,48,48,48,49
  [Character Format]: BPGR0001
  [Node Format]: Cluster ID: unknown Node ID: 1 Node Name: sys2
```

- To display the keys on a VCS failover disk group:

```
# vxfenadm -s /dev/vx/rdmp/disk_8

Reading SCSI Registration Keys...

Device Name: /dev/vx/rdmp/disk_8
Total Number Of Keys: 1
key[0]:
  [Numeric Format]: 65,86,67,83,0,0,0,0
  [Character Format]: AVCS
  [Node Format]: Cluster ID: unknown Node ID: 0 Node Name: sys1
```

- 2 To display the keys that are registered in all the disks specified in a disk file:

```
# vxfenadm -s all -f disk_filename
```

For example:

To display all the keys on coordinator disks:

```
# vxfenadm -s all -f /etc/vxfentab

Device Name: /dev/vx/rdmp/disk_9
Total Number Of Keys: 2
key[0]:
[Numeric Format]: 86,70,70,68,57,52,48,49
[Character Format]: VFFD9401
* [Node Format]: Cluster ID: 64916 Node ID: 1 Node Name: sys2
key[1]:
[Numeric Format]: 86,70,70,68,57,52,48,48
[Character Format]: VFFD9400
* [Node Format]: Cluster ID: 64916 Node ID: 0 Node Name: sys1
```

You can verify the cluster ID using the `lltstat -C` command, and the node ID using the `lltstat -N` command. For example:

```
# lltstat -C
57069
```

If the disk has keys which do not belong to a specific cluster, then the `vxfenadm` command cannot look up the node name for the node ID and hence prints the node name as unknown. For example:

```
Device Name: /dev/vx/rdmp/disk_7
Total Number Of Keys: 1
key[0]:
[Numeric Format]: 86,70,45,45,45,45,48,49
[Character Format]: VF----01
[Node Format]: Cluster ID: unknown Node ID: 1 Node Name: sys2
```

For disks with arbitrary format of keys, the `vxfenadm` command prints all the fields as unknown. For example:

```
[Numeric Format]: 65,66,67,68,49,50,51,45
[Character Format]: ABCD123-
[Node Format]: Cluster ID: unknown Node ID: unknown
Node Name: unknown
```

Verifying that the nodes see the same disk

To confirm whether a disk (or LUN) supports SCSI-3 persistent reservations, two nodes must simultaneously have access to the same disks. Because a shared disk is likely to have a different name on each node, check the serial number to verify the identity of the disk. Use the `vxfenadm` command with the `-i` option to verify that the same serial number for the LUN is returned on all paths to the LUN.

For example, an EMC disk is accessible by the `/dev/vx/rdmp/c2t13d0` path on node A and the `/dev/vx/rdmp/c2t11d0` path on node B.

To verify that the nodes see the same disks

- 1 Verify the connection of the shared storage for data to two of the nodes on which you installed SF Oracle RAC.
- 2 From node A, enter the following command:

```
# vxfenadm -i /dev/vx/rdmp/c2t13d0

Vendor id       : EMC
Product id      : SYMMETRIX
Revision        : 5567
Serial Number   : 42031000a
```

The same serial number information should appear when you enter the equivalent command on node B using the `/dev/vx/rdmp/c2t11d0` path.

On a disk from another manufacturer, Hitachi Data Systems, the output is different and may resemble:

```
# vxfenadm -i /dev/vx/rdmp/c2t1d0

Vendor id       : HITACHI
Product id      : OPEN-3      -HP
Revision        : 0117
Serial Number   : 0401EB6F0002
```

Refer to the `vxfenadm(1M)` manual page for more information.

About the `vxfcntlpre` utility

You can use the `vxfcntlpre` utility to remove SCSI-3 registrations and reservations on the disks.

See [“Removing preexisting keys”](#) on page 140.

This utility currently does not support server-based fencing. You must manually resolve any preexisting split-brain with server-based fencing configuration.

See [“Issues during fencing startup on SF Oracle RAC cluster nodes set up for server-based fencing”](#) on page 211.

Removing preexisting keys

If you encountered a split-brain condition, use the `vxfenclearpre` utility to remove SCSI-3 registrations and reservations on the coordinator disks as well as on the data disks in all shared disk groups.

You can also use this procedure to remove the registration and reservation keys created by another node from a disk.

To clear keys after split-brain

- 1 Stop VCS on all nodes.

```
# hastop -all
```

- 2 Make sure that the port `h` is closed on all the nodes. Run the following command on each node to verify that the port `h` is closed:

```
# gabconfig -a
```

Port `h` must not appear in the output.

- 3 Stop I/O fencing on all nodes. Enter the following command on each node:

```
# /sbin/init.d/vxfen stop
```

- 4 If you have any applications that run outside of VCS control that have access to the shared storage, then shut down all other nodes in the cluster that have access to the shared storage. This prevents data corruption.

- 5 Start the `vxfenclearpre` script:

```
# /opt/VRTSvcs/vxfen/bin/vxfenclearpre
```

6 Read the script's introduction and warning. Then, you can choose to let the script run.

```
Do you still want to continue: [y/n] (default : n) y
```

In some cases, informational messages resembling the following may appear on the console of one of the nodes in the cluster when a node is ejected from a disk/LUN. You can ignore these informational messages.

```
<date> <system name> scsi: WARNING: /sbus@3,0/lpfs@0,0/  
sd@0,1(sd91):  
<date> <system name> Error for Command: <undecoded  
cmd 0x5f> Error Level: Informational  
<date> <system name> scsi: Requested Block: 0 Error Block 0  
<date> <system name> scsi: Vendor: <vendor> Serial Number:  
0400759B006E  
<date> <system name> scsi: Sense Key: Unit Attention  
<date> <system name> scsi: ASC: 0x2a (<vendor unique code  
0x2a>), ASCQ: 0x4, FRU: 0x0
```

The script cleans up the disks and displays the following status messages.

```
Cleaning up the coordinator disks...
```

```
Cleaning up the data disks for all shared disk groups...
```

```
Successfully removed SCSI-3 persistent registration and  
reservations from the coordinator disks as well as the  
shared data disks.
```

```
You can retry starting fencing module. In order to  
restart the whole product, you might want to  
reboot the system.
```

7 Start the fencing module.

```
# /sbin/init.d/vxfen start
```

8 Start VCS on all nodes.

```
# hstart
```

About the vxfsnwap utility

The vxfsnwap utility allows you to add, remove, and replace coordinator points in a cluster that is online. The utility verifies that the serial number of the new disks are identical on all the nodes and the new disks can support I/O fencing.

This utility supports both disk-based and server-based fencing.

Refer to the `vxfsnwap(1M)` manual page.

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for details on the I/O fencing requirements.

You can replace the coordinator disks without stopping I/O fencing in the following cases:

- The disk becomes defective or inoperable and you want to switch to a new diskgroup.

See [“Replacing I/O fencing coordinator disks when the cluster is online”](#) on page 143.

See [“Replacing the coordinator disk group in a cluster that is online”](#) on page 147.

If you want to replace the coordinator disks when the cluster is offline, you cannot use the vxfsnwap utility. You must manually perform the steps that the utility does to replace the coordinator disks.

See [“Replacing defective disks when the cluster is offline”](#) on page 206.

- You want to switch the disk interface between raw devices and DMP devices.

- The keys that are registered on the coordinator disks are lost.

In such a case, the cluster might panic when a network partition occurs. You can replace the coordinator disks with the same disks using the vxfsnwap command. During the disk replacement, the missing keys register again without any risk of data corruption.

See [“Refreshing lost keys on coordinator disks”](#) on page 152.

In server-based fencing configuration, you can use the vxfsnwap utility to perform the following tasks:

- Perform a planned replacement of customized coordination points (CP servers or SCSI-3 disks).

See [“Replacing coordination points for server-based fencing in an online cluster”](#) on page 158.

- Refresh the I/O fencing keys that are registered on the coordination points.

See [“Refreshing registration keys on the coordination points for server-based fencing”](#) on page 156.

You can also use the `vxfsnwap` utility to migrate between the disk-based and the server-based fencing without incurring application downtime in the SF Oracle RAC cluster.

See [“Migrating from disk-based to server-based fencing in an online cluster”](#) on page 60.

See [“Migrating from server-based to disk-based fencing in an online cluster”](#) on page 65.

If the `vxfsnwap` operation is unsuccessful, then you can use the `-a cancel` of the `vxfsnwap` command to manually roll back the changes that the `vxfsnwap` utility does.

- For disk-based fencing, use the `vxfsnwap -g diskgroup -a cancel` command to cancel the `vxfsnwap` operation.
You must run this command if a node fails during the process of disk replacement, or if you aborted the disk replacement.
- For server-based fencing, use the `vxfsnwap -a cancel` command to cancel the `vxfsnwap` operation.

Replacing I/O fencing coordinator disks when the cluster is online

Review the procedures to add, remove, or replace one or more coordinator disks in a cluster that is operational.

Warning: The cluster might panic if any node leaves the cluster membership before the `vxfsnwap` script replaces the set of coordinator disks.

To replace a disk in a coordinator diskgroup when the cluster is online

- 1 Make sure system-to-system communication is functioning properly.
- 2 Determine the value of the `FaultTolerance` attribute.

```
# hares -display coordpoint -attribute FaultTolerance -localclus
```
- 3 Estimate the number of coordination points you plan to use as part of the fencing configuration.

- 4 Set the value of the FaultTolerance attribute to 0.

Note: It is necessary to set the value to 0 because later in the procedure you need to reset the value of this attribute to a value that is lower than the number of coordination points. This ensures that the Coordpoint Agent does not fault.

- 5 Check the existing value of the LevelTwoMonitorFreq attribute.

```
#hares -display coordpoint -attribute LevelTwoMonitorFreq -localclus
```

Note: Make a note of the attribute value before you proceed to the next step. After migration, when you re-enable the attribute you want to set it to the same value.

You can also run the `hares -display coordpoint` to find out whether the LevelTwoMonitorFreq value is set.

- 6 Disable level two monitoring of CoordPoint agent.

```
# hares -modify coordpoint LevelTwoMonitorFreq 0
```

- 7 Make sure that the cluster is online.

```
# vxfenadm -d
```

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: SCSI3
Fencing SCSI3 Disk Policy: dmp
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```


8 Import the coordinator disk group.

The file `/etc/vxfendg` includes the name of the disk group (typically, `vxfencoordg`) that contains the coordinator disks, so use the command:

```
# vxdg -tfc import `cat /etc/vxfendg`
```

where:

-t specifies that the disk group is imported only until the node restarts.

-f specifies that the import is to be done forcibly, which is necessary if one or more disks is not accessible.

-C specifies that any import locks are removed.

9 If your setup uses VRTSvxvm *version*, then skip to step 10. You need not set `coordinator=off` to add or remove disks. For other VxVM versions, perform this step:

Where `<version>` is the specific release version.

Turn off the coordinator attribute value for the coordinator disk group.

```
# vxdg -g vxfencoordg set -o coordinator=off
```

10 To remove disks from the coordinator disk group, use the VxVM disk administrator utility `vxdiskadm`.**11** Perform the following steps to add new disks to the coordinator disk group:

- Add new disks to the node.
- Initialize the new disks as VxVM disks.
- Check the disks for I/O fencing compliance.
- Add the new disks to the coordinator disk group and set the coordinator attribute value as "on" for the coordinator disk group.

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for detailed instructions.

Note that though the disk group content changes, the I/O fencing remains in the same state.

12 From one node, start the `vxfenswap` utility. You must specify the diskgroup to the utility.

The utility performs the following tasks:

- Backs up the existing `/etc/vxfentab` file.

- Creates a test file `/etc/vxfentab.test` for the diskgroup that is modified on each node.
 - Reads the diskgroup you specified in the `vxfenswap` command and adds the diskgroup to the `/etc/vxfentab.test` file on each node.
 - Verifies that the serial number of the new disks are identical on all the nodes. The script terminates if the check fails.
 - Verifies that the new disks can support I/O fencing on each node.
- 13** If the disk verification passes, the utility reports success and asks if you want to commit the new set of coordinator disks.
- 14** Confirm whether you want to clear the keys on the coordination points and proceed with the `vxfenswap` operation.

```
Do you want to clear the keys on the coordination points
and proceed with the vxfenswap operation? [y/n] (default: n) y
```

- 15** Review the message that the utility displays and confirm that you want to commit the new set of coordinator disks. Else skip to step 16.

```
Do you wish to commit this change? [y/n] (default: n) y
```

If the utility successfully commits, the utility moves the `/etc/vxfentab.test` file to the `/etc/vxfentab` file.

- 16** If you do not want to commit the new set of coordinator disks, answer `n`.
The `vxfenswap` utility rolls back the disk replacement operation.
- 17** Re-enable the `LevelTwoMonitorFreq` attribute of the `CoordPoint` agent. You may want to use the value that was set before disabling the attribute.

```
# hares -modify coordpoint LevelTwoMonitorFreq Frequencyvalue
```

where *Frequencyvalue* is the value of the attribute.

- 18** Set the `FaultTolerance` attribute to a value that is lower than 50% of the total number of coordination points.

For example, if there are four (4) coordination points in your configuration, then the attribute value must be lower than two (2). If you set it to a higher value than two (2) the `CoordPoint` agent faults.

Replacing the coordinator disk group in a cluster that is online

You can also replace the coordinator disk group using the `vxfsenswap` utility. The following example replaces the coordinator disk group `vxfencoordg` with a new disk group `vx fendg`.

To replace the coordinator disk group

- 1 Make sure system-to-system communication is functioning properly.
- 2 Determine the value of the `FaultTolerance` attribute.

```
# hares -display coordpoint -attribute FaultTolerance -localclus
```

- 3 Estimate the number of coordination points you plan to use as part of the fencing configuration.
- 4 Set the value of the `FaultTolerance` attribute to 0.

Note: It is necessary to set the value to 0 because later in the procedure you need to reset the value of this attribute to a value that is lower than the number of coordination points. This ensures that the Coordpoint Agent does not fault.

- 5 Check the existing value of the `LevelTwoMonitorFreq` attribute.

```
# hares -display coordpoint -attribute LevelTwoMonitorFreq -localclus
```

Note: Make a note of the attribute value before you proceed to the next step. After migration, when you re-enable the attribute you want to set it to the same value.

- 6 Disable level two monitoring of CoordPoint agent.

```
# hares -modify coordpoint LevelTwoMonitorFreq 0
```

7 Make sure that the cluster is online.

```
# vxfenadm -d
```

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: SCSI3
Fencing SCSI3 Disk Policy: dmp
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

8 Find the name of the current coordinator disk group (typically vxfencoorddg) that is in the /etc/vxfendg file.

```
# cat /etc/vxfendg
vxfencoorddg
```

9 Find the alternative disk groups available to replace the current coordinator disk group.

```
# vxdisk -o alldgs list
```

DEVICE	TYPE	DISK	GROUP	STATUS
c4t0d1	auto:cdsdisk	-	(vxfendg)	online
c4t0d2	auto:cdsdisk	-	(vxfendg)	online
c4t0d3	auto:cdsdisk	-	(vxfendg)	online
c4t0d4	auto:cdsdisk	-	(vxfencoorddg)	online
c4t0d5	auto:cdsdisk	-	(vxfencoorddg)	online
c4t0d6	auto:cdsdisk	-	(vxfencoorddg)	online

10 Validate the new disk group for I/O fencing compliance. Run the following command:

```
# vxfentsthdw -c vxfendg
```

See [“Testing the coordinator disk group using vxfentsthdw -c option”](#) on page 129.

- 11 If the new disk group is not already deported, run the following command to deport the disk group:

```
# vxdg deport vxfendg
```

- 12 Perform one of the following:

- Create the `/etc/vxfenmode.test` file with new fencing mode and disk policy information.
- Edit the existing the `/etc/vxfenmode` with new fencing mode and disk policy information and remove any preexisting `/etc/vxfenmode.test` file.

Note that the format of the `/etc/vxfenmode.test` file and the `/etc/vxfenmode` file is the same.

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for more information.

- 13 From any node, start the `vxfsnwap` utility. For example, if `vxfendg` is the new disk group that you want to use as the coordinator disk group:

```
# vxfsnwap -g vxfendg [-n]
```

The utility performs the following tasks:

- Backs up the existing `/etc/vxfentab` file.
 - Creates a test file `/etc/vxfentab.test` for the disk group that is modified on each node.
 - Reads the disk group you specified in the `vxfsnwap` command and adds the disk group to the `/etc/vxfentab.test` file on each node.
 - Verifies that the serial number of the new disks are identical on all the nodes. The script terminates if the check fails.
 - Verifies that the new disk group can support I/O fencing on each node.
- 14 If the disk verification passes, the utility reports success and asks if you want to replace the coordinator disk group.
 - 15 Confirm whether you want to clear the keys on the coordination points and proceed with the `vxfsnwap` operation.

```
Do you want to clear the keys on the coordination points  
and proceed with the vxfsnwap operation? [y/n] (default: n) y
```

- 16** Review the message that the utility displays and confirm that you want to replace the coordinator disk group. Else skip to step 19.

```
Do you wish to commit this change? [y/n] (default: n) y
```

If the utility successfully commits, the utility moves the `/etc/vxfentab.test` file to the `/etc/vxfentab` file.

The utility also updates the `/etc/vxfendg` file with this new disk group.

- 17** Set the coordinator attribute value as "on" for the new coordinator disk group.

```
# vxdg -g vxfendg set -o coordinator=on
```

Set the coordinator attribute value as "off" for the old disk group.

```
# vxdg -g vxfencoordg set -o coordinator=off
```

- 18** Verify that the coordinator disk group has changed.

```
# cat /etc/vxfendg  
vxfendg
```

The swap operation for the coordinator disk group is complete now.

- 19** If you do not want to replace the coordinator disk group, answer n at the prompt.

The `vxfenswap` utility rolls back any changes to the coordinator disk group.

- 20** Re-enable the `LevelTwoMonitorFreq` attribute of the `CoordPoint` agent. You may want to use the value that was set before disabling the attribute.

```
# hares -modify coordpoint LevelTwoMonitorFreq Frequencyvalue
```

where *Frequencyvalue* is the value of the attribute.

- 21** Set the `FaultTolerance` attribute to a value that is lower than 50% of the total number of coordination points.

For example, if there are four (4) coordination points in your configuration, then the attribute value must be lower than two (2). If you set it to a higher value than two (2) the `CoordPoint` agent faults.

Adding disks from a recovered site to the coordinator diskgroup

In a campus cluster environment, consider a case where the primary site goes down and the secondary site comes online with a limited set of disks. When the primary site restores, the primary site's disks are also available to act as

coordinator disks. You can use the `vxferswap` utility to add these disks to the coordinator diskgroup.

To add new disks from a recovered site to the coordinator diskgroup

- 1 Make sure system-to-system communication is functioning properly.
- 2 Make sure that the cluster is online.

```
# vxferadm -d
```

```
I/O Fencing Cluster Information:
```

```
=====
```

```
Fencing Protocol Version: 201
```

```
Fencing Mode: SCSI3
```

```
Fencing SCSI3 Disk Policy: dmp
```

```
Cluster Members:
```

```
  * 0 (sys1)
```

```
  1 (sys2)
```

```
RFSM State Information:
```

```
  node 0 in state 8 (running)
```

```
  node 1 in state 8 (running)
```

- 3 Verify the name of the coordinator diskgroup.

```
# cat /etc/vxfendg
```

```
vxfercoorddg
```

- 4 Run the following command:

```
# vxferdisk -o alldgs list
```

DEVICE	TYPE	DISK	GROUP	STATUS
c1t1d0	auto:cdsdisk	-	(vxfercoorddg)	online
c2t1d0	auto	- -	offline	
c3t1d0	auto	- -	offline	

5 Verify the number of disks used in the coordinator diskgroup.

```
# vxfenconfig -l
I/O Fencing Configuration Information:
=====
Count                : 1
Disk List
Disk Name            Major  Minor  Serial Number      Policy
/dev/vx/rdmp/c1t1d0      32   48   R450 00013154 0312      dmp
```

6 When the primary site comes online, start the vxfenswap utility on any node in the cluster:

```
# vxfenswap -g vxfencoordg [-n]
```

7 Verify the count of the coordinator disks.

```
# vxfenconfig -l
I/O Fencing Configuration Information:
=====
Single Disk Flag      : 0
Count                 : 3
Disk List
Disk Name            Major  Minor  Serial Number      Policy
/dev/vx/rdmp/c1t1d0      32   48   R450 00013154 0312      dmp
/dev/vx/rdmp/c2t1d0      32   32   R450 00013154 0313      dmp
/dev/vx/rdmp/c3t1d0      32   16   R450 00013154 0314      dmp
```

Refreshing lost keys on coordinator disks

If the coordinator disks lose the keys that are registered, the cluster might panic when a network partition occurs.

You can use the vxfenswap utility to replace the coordinator disks with the same disks. The vxfenswap utility registers the missing keys during the disk replacement.

To refresh lost keys on coordinator disks

- 1 Make sure system-to-system communication is functioning properly.
- 2 Make sure that the cluster is online.

```
# vxfenadm -d
```

```
I/O Fencing Cluster Information:
=====
Fencing Protocol Version: 201
Fencing Mode: SCSI3
Fencing SCSI3 Disk Policy: dmp
Cluster Members:
  * 0 (sys1)
  1 (sys2)
RFSM State Information:
  node 0 in state 8 (running)
  node 1 in state 8 (running)
```

- 3 Run the following command to view the coordinator disks that do not have keys:

```
# vxfenadm -s all -f /etc/vxfentab
```

```
Device Name: /dev/vx/rdmp/clt1d0
Total Number of Keys: 0
No keys...
...
```

- 4 Copy the `/etc/vxfenmode` file to the `/etc/vxfenmode.test` file.

This ensures that the configuration details of both the files are the same.

- 5 On any node, run the following command to start the `vxfenswap` utility:

```
# vxfenswap -g vxfencoorddg [-n]
```

- 6 Verify that the keys are atomically placed on the coordinator disks.

```
# vxfenadm -s all -f /etc/vxfentab
```

```
Device Name: /dev/vx/rdmp/clt1d0
```

```
Total Number of Keys: 4
```

```
...
```

Enabling or disabling the preferred fencing policy

You can enable or disable the preferred fencing feature for your I/O fencing configuration.

You can enable preferred fencing to use system-based race policy or group-based race policy. If you disable preferred fencing, the I/O fencing configuration uses the default count-based race policy.

See [“About preferred fencing”](#) on page 44.

To enable preferred fencing for the I/O fencing configuration

- 1 Make sure that the cluster is running with I/O fencing set up.

```
# vxfenadm -d
```

- 2 Make sure that the cluster-level attribute `UseFence` has the value set to `SCSI3`.

```
# haclus -value UseFence
```

- 3 To enable system-based race policy, perform the following steps:

- Make the VCS configuration writable.

```
# haconf -makerw
```

- Set the value of the cluster-level attribute `PreferredFencingPolicy` as `System`.

```
# haclus -modify PreferredFencingPolicy System
```

- Set the value of the system-level attribute `FencingWeight` for each node in the cluster.

For example, in a two-node cluster, where you want to assign `sys1` five times more weight compared to `sys2`, run the following commands:

```
# hasys -modify sys1 FencingWeight 50
# hasys -modify sys2 FencingWeight 10
```

- Save the VCS configuration.

```
# haconf -dump -makero
```

- 4 To enable group-based race policy, perform the following steps:

- Make the VCS configuration writable.

```
# haconf -makerw
```

- Set the value of the cluster-level attribute `PreferredFencingPolicy` as `Group`.

```
# haclus -modify PreferredFencingPolicy Group
```

- Set the value of the group-level attribute `Priority` for each service group. For example, run the following command:

```
# hagrps -modify service_group Priority 1
```

Make sure that you assign a parent service group an equal or lower priority than its child service group. In case the parent and the child service groups are hosted in different subclusters, then the subcluster that hosts the child service group gets higher preference.

- Save the VCS configuration.

```
# haconf -dump -makero
```

- 5 To view the fencing node weights that are currently set in the fencing driver, run the following command:

```
# vxfenconfig -a
```

To disable preferred fencing for the I/O fencing configuration

- 1 Make sure that the cluster is running with I/O fencing set up.

```
# vxfenadm -d
```

- 2 Make sure that the cluster-level attribute UseFence has the value set to SCSI3.

```
# haclus -value UseFence
```

- 3 To disable preferred fencing and use the default race policy, set the value of the cluster-level attribute PreferredFencingPolicy as Disabled.

```
# haconf -makerw
```

```
# haclus -modify PreferredFencingPolicy Disabled
```

```
# haconf -dump -makero
```

Administering the CP server

This section provides the following CP server administration information:

- CP server administration user types and privileges
- CP server administration command (cpsadm)

This section also provides instructions for the following CP server administration tasks:

- Refreshing registration keys on the coordination points for server-based fencing
- Coordination Point replacement for an online cluster
- Migrating from non-secure to secure setup for CP server and SF Oracle RAC cluster communication

Refreshing registration keys on the coordination points for server-based fencing

Replacing keys on a coordination point (CP server) when the SF Oracle RAC cluster is online involves refreshing that coordination point's registrations. You can perform a planned refresh of registrations on a CP server without incurring application downtime on the SF Oracle RAC cluster. You must refresh registrations on a CP server if the CP server agent issues an alert on the loss of such registrations on the CP server database.

The following procedure describes how to refresh the coordination point registrations.

To refresh the registration keys on the coordination points for server-based fencing

- 1 Ensure that the SF Oracle RAC cluster nodes and users have been added to the new CP server(s). Run the following commands:

```
# cpsadm -s cp_server -a list_nodes  
  
# cpsadm -s cp_server -a list_users
```

If the SF Oracle RAC cluster nodes are not present here, prepare the new CP server(s) for use by the SF Oracle RAC cluster.

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for instructions.

- 2 Ensure that fencing is running on the cluster in customized mode using the coordination points mentioned in the `/etc/vxfenmode` file.

If the `/etc/vxfenmode.test` file exists, ensure that the information in it and the `/etc/vxfenmode` file are the same. Otherwise, `vxfenswap` utility uses information listed in `/etc/vxfenmode.test` file.

For example, enter the following command:

```
# vxfenadm -d  
  
=====
```

```
Fencing Protocol Version: 201  
Fencing Mode: CUSTOMIZED  
Cluster Members:  
* 0 (galaxy)  
1 (nebula)  
RFSM State Information:  
node 0 in state 8 (running)  
node 1 in state 8 (running)
```

- 3 List the coordination points currently used by I/O fencing :

```
# vxfenconfig -l
```

- 4 Copy the `/etc/vxfenmode` file to the `/etc/vxfenmode.test` file.

This ensures that the configuration details of both the files are the same.

5 Run the `vxfsnwap` utility from one of the nodes of the cluster.

The `vxfsnwap` utility requires secure ssh connection to all the cluster nodes. Use `-n` to use rsh instead of default ssh.

For example:

```
# vxfsnwap [-n]
```

The command returns:

```
VERITAS vxfsnwap version <version> <platform>
The logfile generated for vxfsnwap is
/var/VRTSvcs/log/vxfen/vxfsnwap.log.
19156
Please Wait...
VXFEN vxfsnconfig NOTICE Driver will use customized fencing
- mechanism cps
Validation of coordination points change has succeeded on
all nodes.
You may commit the changes now.
WARNING: This may cause the whole cluster to panic
if a node leaves membership before the change is complete.
```

6 You are then prompted to commit the change. Enter `y` for yes.

The command returns a confirmation of successful coordination point replacement.

7 Confirm the successful execution of the `vxfsnwap` utility. If CP agent is configured, it should report ONLINE as it succeeds to find the registrations on coordination points. The registrations on the CP server and coordinator disks can be viewed using the `cpsadm` and `vxfsnadm` utilities respectively.

Note that a running online coordination point refreshment operation can be canceled at any time using the command:

```
# vxfsnwap -a cancel
```

Replacing coordination points for server-based fencing in an online cluster

Use the following procedure to perform a planned replacement of customized coordination points (CP servers or SCSI-3 disks) without incurring application downtime on an online SF Oracle RAC cluster.

Note: If multiple clusters share the same CP server, you must perform this replacement procedure in each cluster.

You can use the `vxfenswap` utility to replace coordination points when fencing is running in customized mode in an online cluster, with `vxfen_mechanism=cps`. The utility also supports migration from server-based fencing (`vxfen_mode=customized`) to disk-based fencing (`vxfen_mode=scsi3`) and vice versa in an online cluster.

However, if the SF Oracle RAC cluster has fencing disabled (`vxfen_mode=disabled`), then you must take the cluster offline to configure disk-based or server-based fencing.

See [“Deployment and migration scenarios for CP server”](#) on page 55.

You can cancel the coordination point replacement operation at any time using the `vxfenswap -a cancel` command.

See [“About the vxfenswap utility”](#) on page 142.

To replace coordination points for an online cluster

- 1 Ensure that the SF Oracle RAC cluster nodes and users have been added to the new CP server(s). Run the following commands:

```
# cpsadm -s cpserver -a list_nodes  
# cpsadm -s cpserver -a list_users
```

If the SF Oracle RAC cluster nodes are not present here, prepare the new CP server(s) for use by the SF Oracle RAC cluster.

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for instructions.

- 2 Ensure that fencing is running on the cluster using the old set of coordination points and in customized mode.

For example, enter the following command:

```
# vxfenadm -d
```

The command returns:

```
I/O Fencing Cluster Information:  
=====  
Fencing Protocol Version: <version>  
Fencing Mode: Customized  
Cluster Members:  
* 0 (sys1)  
1 (sys2)  
RFSM State Information:  
node 0 in state 8 (running)  
node 1 in state 8 (running)
```

- 3 Create a new `/etc/vxfenmode.test` file on each SF Oracle RAC cluster node with the fencing configuration changes such as the CP server information.

Review and if necessary, update the `vxfenmode` parameters for security, the coordination points, and if applicable to your configuration, `vxfendg`.

Refer to the text information within the `vxfenmode` file for additional information about these parameters and their new possible values.

- 4 From one of the nodes of the cluster, run the `vxfenswap` utility.

The `vxfenswap` utility requires secure ssh connection to all the cluster nodes. Use `-n` to use rsh instead of default ssh.

```
# vxfenswap [-n]
```


- 5 Review the message that the utility displays and confirm whether you want to commit the change.

- If you do not want to commit the new fencing configuration changes, press Enter or answer n at the prompt.

```
Do you wish to commit this change? [y/n] (default: n) n
```

The `vxferswap` utility rolls back the migration operation.

- If you want to commit the new fencing configuration changes, answer y at the prompt.

```
Do you wish to commit this change? [y/n] (default: n) y
```

If the utility successfully completes the operation, the utility moves the `/etc/vxfermode.test` file to the `/etc/vxfermode` file.

- 6 Confirm the successful execution of the `vxferswap` utility by checking the coordination points currently used by the `vxfer` driver.

For example, run the following command:

```
# vxferconfig -l
```

Migrating from non-secure to secure setup for CP server and SF Oracle RAC cluster communication

The following procedure describes how to migrate from a non-secure to secure set up for the coordination point server (CP server) and SF Oracle RAC cluster.

To migrate from non-secure to secure setup for CP server and SF Oracle RAC cluster

- 1 Stop VCS on all cluster nodes that use the CP servers.

```
# hastop -all
```

- 2 Stop fencing on all the SF Oracle RAC cluster nodes of all the clusters.

```
# /sbin/init.d/vxfer stop
```

- 3 Stop all the CP servers using the following command on each CP server:

```
# hagrps -offline CPSSG -any
```

- 4 Ensure that security is configured for communication on CP Servers as well as all the clients.

See the *Veritas Storage Foundation for Oracle RAC Installation Guide* for more information.

- 5 ■ If CP server is hosted on an SFHA cluster, perform this step on each CP server.
Bring the mount resource in the CPSSG service group online.

```
# hares -online cpsmount -sys local_system_name
```

Complete the remaining steps.

- If CP server is hosted on a single-node VCS cluster, skip to step 8 and complete the remaining steps.

- 6 After the mount resource comes online, move the `credentials` directory from the default location to shared storage.

```
# mv /var/VRTSvcs/vcsauth/data/CPSEVER /etc/VRTSvcs/db/
```

- 7 Create softlinks on all the nodes of the CP servers.

```
# ln -s /etc/VRTScps/db/CPSEVER \  
/var/VRTSvcs/vcsauth/data/CPSEVER
```

- 8 Edit `/etc/vxcps.conf` on each CP server to set `security=1`.

- 9 Start CP servers by using the following command:

```
# hagrps -online CPSSG -any
```

- 10 Edit `/etc/VRTSvcs/conf/config/main.cf` on the first node of the cluster and remove the `UseFence=SCSI3` attribute.

Start VCS on the first node and then on all other nodes of the cluster.

- 11 Reconfigure fencing on each cluster by using the installer.

```
# /opt/VRTS/install/installsfrac<version> -fencing
```

Where `<version>` is the specific release version.

Administering CFS

This section describes some of the major aspects of Cluster File System (CFS) administration.

This section provides instructions for the following CFS administration tasks:

- Resizing CFS file systems
See [“Resizing CFS file systems”](#) on page 163.
- Verifying the status of CFS file systems
See [“Verifying the status of CFS file system nodes and their mount points”](#) on page 163.

If you encounter issues while administering CFS, refer to the troubleshooting section for assistance.

Resizing CFS file systems

If you see a message on the console indicating that a Cluster File System (CFS) file system is full, you may want to resize the file system. The `vxresize` command lets you resize a CFS file system. It extends the file system and the underlying volume.

See the `vxresize (1M)` manual page for information on various options.

The following command resizes an Oracle data CFS file system (the Oracle data volume is CFS mounted):

```
# vxresize -g oradatadg oradatavol +2G
```

where `oradatadg` is the CVM disk group, `oradatavol` is the volume, and `+2G` indicates the increase in volume size by 2 Gigabytes.

Verifying the status of CFS file system nodes and their mount points

Run the `cfsccluster status` command to see the status of the nodes and their mount points:

```
# cfsccluster status
```

```
Node           : sys1
Cluster Manager : running
CVM state      : running
```

MOUNT POINT	SHARED VOLUME	DISK GROUP	STATUS
/app/crshome	crsbinvol	bindg	MOUNTED
/ocrvote	ocrvotevol	ocrvotedg	MOUNTED
/app/oracle/orahome	orabinvol	bindg	MOUNTED
/oradata1	oradatavol	oradatadg	MOUNTED
/arch	archvol	oradatadg	MOUNTED

```
Node           : sys2
Cluster Manager : running
CVM state      : running
```

MOUNT POINT	SHARED VOLUME	DISK GROUP	STATUS
/app/crshome	crsbinvol	bindg	MOUNTED
/ocrvote	ocrvotevol	ocrvotedg	MOUNTED
/app/oracle/orahome	orabinvol	bindg	MOUNTED
/oradata1	oradatavol	oradatadg	MOUNTED
/arch	archvol	oradatadg	MOUNTED

Administering CVM

This section provides instructions for the following CVM administration tasks:

- Establishing CVM cluster membership manually
See [“Establishing CVM cluster membership manually”](#) on page 165.
- Changing CVM master manually
See [“Changing the CVM master manually”](#) on page 165.
- Importing a shared disk group manually
See [“Importing a shared disk group manually”](#) on page 168.
- Deporting a shared disk group manually
See [“Deporting a shared disk group manually”](#) on page 169.
- Verifying if CVM is running in an SF Oracle RAC cluster
See [“Verifying if CVM is running in an SF Oracle RAC cluster”](#) on page 169.
- Verifying CVM membership state
See [“Verifying CVM membership state”](#) on page 170.
- Verifying the state of CVM shared disk groups
See [“Verifying the state of CVM shared disk groups”](#) on page 170.
- Verifying the activation mode
See [“Verifying the activation mode”](#) on page 170.

If you encounter issues while administering CVM, refer to the troubleshooting section for assistance.

See [“Troubleshooting Cluster Volume Manager in SF Oracle RAC clusters”](#) on page 213.

Listing all the CVM shared disks

You can use the following command to list all the CVM shared disks:

```
# vxdisk -o alldgs list |grep shared
```

Establishing CVM cluster membership manually

In most cases you do not have to start CVM manually; it normally starts when VCS is started.

Run the following command to start CVM manually:

```
# vxclustadm -m vcs -t gab startnode
```

```
vxclustadm: initialization completed
```

Note that `vxclustadm` reads `main.cf` for cluster configuration information and is therefore not dependent upon VCS to be running. You do not need to run the `vxclustadm startnode` command as normally the `hastart` (VCS start) command starts CVM automatically.

To verify whether CVM is started properly:

```
# vxclustadm nidmap
```

Name	CVM Nid	CM Nid	State
sys1	0	0	Joined: Master
sys2	1	1	Joined: Slave

Changing the CVM master manually

You can change the CVM master manually from one node in the cluster to another node, while the cluster is online. CVM migrates the master node, and reconfigures the cluster.

Symantec recommends that you switch the master when the cluster is not handling VxVM configuration changes or cluster reconfiguration operations. In most cases, CVM aborts the operation to change the master, if CVM detects that any configuration changes are occurring in the VxVM or the cluster. After the master change operation starts reconfiguring the cluster, other commands that require configuration changes will fail until the master switch completes.

See [“Errors during CVM master switching”](#) on page 168.

To change the master online, the cluster must be cluster protocol version 100 or greater.

To change the CVM master manually

- 1 To view the current master, use one of the following commands:

```
# vxclustadm nidmap
Name           CVM Nid    CM Nid     State
sys1           0          0          Joined: Slave
sys2           1          1          Joined: Master

# vxdctl -c mode
mode: enabled: cluster active - MASTER
master: sys2
```

In this example, the CVM master is sys2.

- 2 From any node on the cluster, run the following command to change the CVM master:

```
# vxclustadm setmaster nodename
```

where *nodename* specifies the name of the new CVM master.

The following example shows changing the master on a cluster from sys2 to sys1:

```
# vxclustadm setmaster sys1
```

3 To monitor the master switching, use the following command:

```
# vxclustadm -v nodestate
state: cluster member
      nodeId=0
      masterId=0
      neighborId=1
      members[0]=0xf
      joiners[0]=0x0
      leavers[0]=0x0
      members[1]=0x0
      joiners[1]=0x0
      leavers[1]=0x0
      reconfig_seqnum=0x9f9767
      vxfen=off
state: master switching in progress
reconfig: vxconfigd in join
```

In this example, the state indicates that master is being changed.

4 To verify whether the master has successfully changed, use one of the following commands:

```
# vxclustadm nidmap
Name           CVM Nid   CM Nid   State
sys1           0         0        Joined: Master
sys2           1         1        Joined: Slave

# vxdctl -c mode
mode: enabled: cluster active - MASTER
master: sys1
```

Considerations for changing the master manually

If the master is not running on the node best suited to be the master of the cluster, you can manually change the master. Here are some scenarios when this might occur.

- The currently running master lost access to some of its disks.
By default, CVM uses I/O shipping to handle this scenario. However, you may want to failover the application to a node which has access to the disks. When you move the application, you may also want to relocate the master role to a new node. For example, you may want the master node and the application to be on the same node.

You can use the master switching operation to move the master role without causing the original master node to leave the cluster. After the master role and the application are both switched to other nodes, you may want to remove the original node from the cluster. You can unmount the file systems and cleanly shut down the node. You can then do maintenance on the node.

- The master node is not scaling well with the overlap of application load and the internally-generated administrative I/Os.
You may choose to reevaluate the placement strategy and relocate the master node.

Errors during CVM master switching

Symantec recommends that you switch the master when the cluster is not handling VxVM or cluster configuration changes.

In most cases, CVM aborts the operation to change the master, if CVM detects any configuration changes in progress. CVM logs the reason for the failure into the system logs. In some cases, the failure is displayed in the `vxclustadm setmaster` output as follows:

```
# vxclustadm setmaster sys1
VxVM vxclustadm ERROR V-5-1-15837 Master switching, a reconfiguration or
a transaction is in progress.
Try again
```

In some cases, if the master switching operation is interrupted with another reconfiguration operation, the master change fails. In this case, the existing master remains the master of the cluster. After the reconfiguration is complete, reissue the `vxclustadm setmaster` command to change the master.

If the master switching operation has started the reconfiguration, any command that initiates a configuration change fails with the following error:

```
Node processing a master-switch request. Retry operation.
```

If you see this message, retry the command after the master switching has completed.

Importing a shared disk group manually

You can use the following command to manually import a shared disk group:

```
# vxdg -s import dg_name
```


Deporting a shared disk group manually

You can use the following command to manually deport a shared disk group:

```
# vxdg deport dg_name
```

Note that the deport of a shared disk group removes the SCSI-3 PGR keys on the disks.

Starting shared volumes manually

Following a manual CVM shared disk group import, the volumes in the disk group need to be started manually, as follows:

```
# vxvol -g dg_name startall
```

To verify that the volumes are started, run the following command:

```
# vxprint -htrg dg_name | grep ^v
```

Verifying if CVM is running in an SF Oracle RAC cluster

You can use the following options to verify whether CVM is up or not in an SF Oracle RAC cluster.

The following output is displayed on a node that is not a member of the cluster:

```
# vxdctl -c mode
mode: enabled: cluster inactive
# vxclustadm -v nodestate
state: out of cluster
```

On the master node, the following output is displayed:

```
# vxdctl -c mode
mode: enabled: cluster active - MASTER
master: sys1
```

On the slave nodes, the following output is displayed:

```
# vxdctl -c mode
mode: enabled: cluster active - SLAVE
master: sys2
```

The following command lets you view all the CVM nodes at the same time:

```
# vxclustadm nidmap

Name      CVM Nid   CM Nid   State
sys1      0         0        Joined: Master
sys2      1         1        Joined: Slave
```

Verifying CVM membership state

The state of CVM can be verified as follows:

```
# vxclustadm -v nodestate

state: joining
      nodeId=0
      masterId=0
      neighborId=0
      members=0x1
      joiners=0x0
      leavers=0x0
      reconfig_seqnum=0x0
      reconfig: vxconfigd in join
```

The state indicates that CVM has completed its kernel level join and is in the middle of vxconfigd level join.

The `vxctl -c mode` command indicates whether a node is a CVM master or CVM slave.

Verifying the state of CVM shared disk groups

You can use the following command to list the shared disk groups currently imported in the SF Oracle RAC cluster:

```
# vxdg list |grep shared

orabinvol_dg enabled,shared,cds 1052685125.1485.csha3
```

Verifying the activation mode

In an SF Oracle RAC cluster, the activation of shared disk group should be set to “shared-write” on each of the cluster nodes.

To verify whether the “shared-write” activation is set:

```
# vxdg list dg_name |grep activation

local-activation: shared-write
```

If "shared-write" activation is not set, run the following command:

```
# vxdg -g dg_name set activation=sw
```

Administering SF Oracle RAC global clusters

This section provides instructions for the following global cluster administration tasks:

- About setting up a fire drill
See [“About setting up a disaster recovery fire drill”](#) on page 171.
- Configuring the fire drill service group using the wizard
See [“About configuring the fire drill service group using the Fire Drill Setup wizard”](#) on page 172.
- Verifying a successful fire drill
See [“Verifying a successful fire drill”](#) on page 174.
- Scheduling a fire drill
See [“Scheduling a fire drill”](#) on page 174.

For a sample fire drill service group configuration:

See [“Sample fire drill service group configuration”](#) on page 174.

About setting up a disaster recovery fire drill

The disaster recovery fire drill procedure tests the fault-readiness of a configuration by mimicking a failover from the primary site to the secondary site. This procedure is done without stopping the application at the primary site and disrupting user access, interrupting the flow of replicated data, or causing the secondary site to need resynchronization.

The initial steps to create a fire drill service group on the secondary site that closely follows the configuration of the original application service group and contains a point-in-time copy of the production data in the Replicated Volume Group (RVG). Bringing the fire drill service group online on the secondary site demonstrates the ability of the application service group to fail over and come online at the secondary site, should the need arise. Fire drill service groups do not interact with outside clients or with other instances of resources, so they can safely come online even when the application service group is online.

You must conduct a fire drill only at the secondary site; do not bring the fire drill service group online on the node hosting the original application.

Before you perform a fire drill in a disaster recovery setup that uses VVR, perform the following steps:

- Set the value of the ReuseMntPt attribute to 1 for all Mount resources.
- Configure the fire drill service group.
See [“About configuring the fire drill service group using the Fire Drill Setup wizard”](#) on page 172.
- After the fire drill service group is taken offline, reset the value of the ReuseMntPt attribute to 0 for all Mount resources.

VCS also supports HA fire drills to verify a resource can fail over to another node in the cluster.

Note: You can conduct fire drills only on regular VxVM volumes; volume sets (vset) are not supported.

VCS provides hardware replication agents for array-based solutions, such as Hitachi Truecopy, EMC SRDF, and so on . If you are using hardware replication agents to monitor the replicated data clusters, refer to the VCS replication agent documentation for details on setting up and configuring fire drill.

About configuring the fire drill service group using the Fire Drill Setup wizard

Use the Fire Drill Setup Wizard to set up the fire drill configuration.

The wizard performs the following specific tasks:

- Creates a Cache object to store changed blocks during the fire drill, which minimizes disk space and disk spindles required to perform the fire drill.
- Configures a VCS service group that resembles the real application group.

The wizard works only with application groups that contain one disk group. The wizard sets up the first RVG in an application. If the application has more than one RVG, you must create space-optimized snapshots and configure VCS manually, using the first RVG as reference.

You can schedule the fire drill for the service group using the `fdsched` script.

See [“Scheduling a fire drill”](#) on page 174.

Running the fire drill setup wizard

To run the wizard

- 1 Start the RVG Secondary Fire Drill wizard on the VVR secondary site, where the application service group is offline and the replication group is online as a secondary:

```
# /opt/VRTSvcs/bin/fdsetup
```

- 2 Read the information on the Welcome screen and press the **Enter** key.
- 3 The wizard identifies the global service groups. Enter the name of the service group for the fire drill.
- 4 Review the list of volumes in disk group that could be used for a space-optimized snapshot. Enter the volumes to be selected for the snapshot. Typically, all volumes used by the application, whether replicated or not, should be prepared, otherwise a snapshot might not succeed.

Press the **Enter** key when prompted.

- 5 Enter the cache size to store writes when the snapshot exists. The size of the cache must be large enough to store the expected number of changed blocks during the fire drill. However, the cache is configured to grow automatically if it fills up. Enter disks on which to create the cache.

Press the **Enter** key when prompted.

- 6 The wizard starts running commands to create the fire drill setup.

Press the **Enter** key when prompted.

The wizard creates the application group with its associated resources. It also creates a fire drill group with resources for the application (Oracle, for example), the CFMount, and the RVGSnapshot types.

The application resources in both service groups define the same application, the same database in this example. The wizard sets the FireDrill attribute for the application resource to 1 to prevent the agent from reporting a concurrency violation when the actual application instance and the fire drill service group are online at the same time.

About configuring local attributes in the fire drill service group

The fire drill setup wizard does not recognize localized attribute values for resources. If the application service group has resources with local (per-system) attribute values, you must manually set these attributes after running the wizard.

Verifying a successful fire drill

Bring the fire drill service group online on a node that does not have the application running. Verify that the fire drill service group comes online. This action validates that your disaster recovery solution is configured correctly and the production service group will fail over to the secondary site in the event of an actual failure (disaster) at the primary site.

If the fire drill service group does not come online, review the VCS engine log to troubleshoot the issues so that corrective action can be taken as necessary in the production service group.

You can also view the fire drill log, located at `/tmp/fd-servicegroup.pid`

Remember to take the fire drill offline once its functioning has been validated. Failing to take the fire drill offline could cause failures in your environment. For example, if the application service group were to fail over to the node hosting the fire drill service group, there would be resource conflicts, resulting in both service groups faulting.

Scheduling a fire drill

You can schedule the fire drill for the service group using the `fdsched` script. The `fdsched` script is designed to run only on the lowest numbered node that is currently running in the cluster. The scheduler runs the command `hagrp -online firedrill_group -any` at periodic intervals.

To schedule a fire drill

- 1 Add the file `/opt/VRTSvcs/bin/fdsched` to your crontab.
- 2 To make fire drills highly available, add the `fdsched` file to each node in the cluster.

Sample fire drill service group configuration

The sample configuration in this section describes a fire drill service group configuration on the secondary site. The configuration uses VVR for replicating data between the sites.

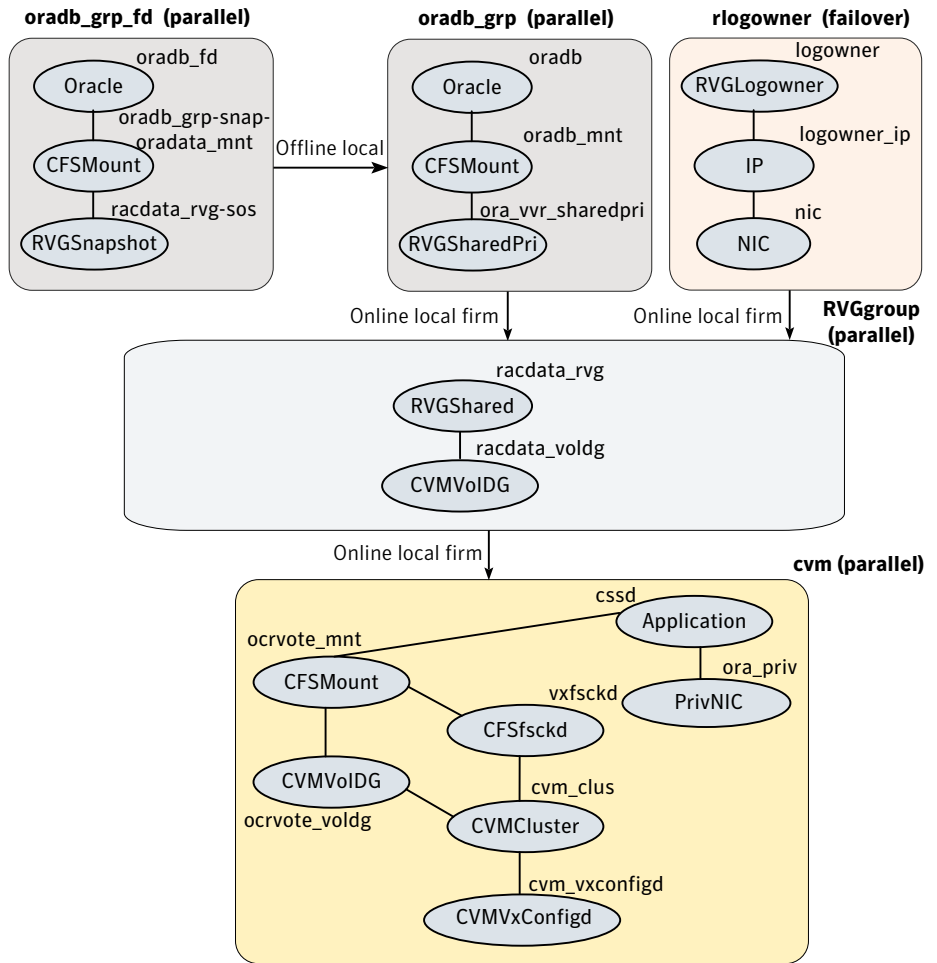
The sample service group describes the following configuration:

- Two SF Oracle RAC clusters, comprising two nodes each, hosted at different geographical locations.
- A single Oracle database that is stored on CFS.
- The database is managed by the VCS agent for Oracle. The agent starts, stops, and monitors the database.

- The database uses the Oracle UDP IPC for database cache fusion.
- A common IP address is used by Oracle Clusterware and database cache fusion. The private IP address is managed by the PrivNIC agent for high availability.
- One virtual IP address must be configured under the `ClusterService` group on each site for inter-cluster communication.
- The Oracle Cluster Registry (OCR) and voting disk are stored on CFS.
- Veritas Volume Replicator (VVR) is used to replicate data between the sites.
- The shared volumes replicated across the sites are configured under the RVG group.
- The replication link used by VVR for communicating log information between sites are configured under the `rlogowner` group. This is a failover group that will be online on only one of the nodes in the cluster at each site.
- The database group is configured as a global group by specifying the clusters on the primary and secondary sites as values for the `ClusterList` group attribute.
- The fire drill service group `oradb_grp_fd` creates a snapshot of the replicated data on the secondary site and starts the database using the snapshot. An offline local dependency is set between the fire drill service group and the application service group to make sure a fire drill does not block an application failover in case a disaster strikes the primary site.

Figure 2-1 illustrates the configuration.

Figure 2-1 Service group configuration for fire drill



Performance and troubleshooting

- [Chapter 3. Troubleshooting SF Oracle RAC](#)
- [Chapter 4. Prevention and recovery strategies](#)
- [Chapter 5. Tunable parameters](#)

Troubleshooting SF Oracle RAC

This chapter includes the following topics:

- [About troubleshooting SF Oracle RAC](#)
- [What to do if you see a licensing reminder](#)
- [Restarting the installer after a failed connection](#)
- [Installer cannot create UUID for the cluster](#)
- [Troubleshooting SF Oracle RAC pre-installation check failures](#)
- [Troubleshooting LLT health check warning messages](#)
- [Troubleshooting LMX health check warning messages in SF Oracle RAC clusters](#)
- [Troubleshooting I/O fencing](#)
- [Troubleshooting Cluster Volume Manager in SF Oracle RAC clusters](#)
- [Troubleshooting VCSIPC](#)
- [Troubleshooting Oracle](#)
- [Troubleshooting ODM in SF Oracle RAC clusters](#)

About troubleshooting SF Oracle RAC

Use the information in this chapter to diagnose setup or configuration problems that you might encounter. For issues that arise from the component products, it may be necessary to refer to the appropriate documentation to resolve it.

Gathering information from an SF Oracle RAC cluster for support analysis

Use troubleshooting scripts to gather information about the configuration and status of your cluster and its modules. The scripts identify depot information, debugging messages, console messages, and information about disk groups and volumes. Forwarding the output of these scripts to Symantec Tech Support can assist with analyzing and solving any problems.

- Gathering configuration information using SORT Data Collector
See “[Gathering configuration information using SORT Data Collector](#)” on page 180.
- Gathering SF Oracle RAC information for support analysis
See “[Gathering SF Oracle RAC information for support analysis](#)” on page 180.
- Gathering VCS information for support analysis
See “[Gathering VCS information for support analysis](#)” on page 181.
- Gathering LLT and GAB information for support analysis
See “[Gathering LLT and GAB information for support analysis](#)” on page 181.
- Gathering IMF information for support analysis
See “[Gathering IMF information for support analysis](#)” on page 182.

Gathering configuration information using SORT Data Collector

SORT Data Collector now supersedes the VRTExplorer utility.

Run the Data Collector with the `VxExplorer` option to gather system and configuration information from a node to diagnose or analyze issues in the cluster.

If you find issues in the cluster that require professional help, run the Data Collector and send the tar file output to Symantec Technical Support to resolve the issue.

Visit the SORT Website and download the UNIX Data Collector appropriate for your operating system:

<https://sort.symantec.com>

For more information:

<https://sort.symantec.com/public/help/wwhelp/wwhimpl/js/html/wwhelp.htm>

Gathering SF Oracle RAC information for support analysis

You must run the `getdbac` script to gather information about the SF Oracle RAC modules.

To gather SF Oracle RAC information for support analysis

- ◆ Run the following command on each node:

```
# /opt/VRTSvcs/bin/getdbac -local
```

The script saves the output to the default file

```
/tmp/vcsopslog.time_stamp.tar.Z
```

If you are unable to resolve the issue, contact Symantec Technical Support with the file.

Gathering VCS information for support analysis

You must run the `hagetcf` command to gather information when you encounter issues with VCS. Symantec Technical Support uses the output of these scripts to assist with analyzing and solving any VCS problems. The `hagetcf` command gathers information about the installed software, cluster configuration, systems, logs, and related information and creates a gzip file.

See the `hagetcf(1M)` manual page for more information.

To gather VCS information for support analysis

- ◆ Run the following command on each node:

```
# /opt/VRTSvcs/bin/hagetcf
```

The command prompts you to specify an output directory for the gzip file. You may save the gzip file to either the default `/tmp` directory or a different directory.

Troubleshoot and fix the issue.

If the issue cannot be fixed, then contact Symantec technical support with the file that the `hagetcf` command generates.

Gathering LLT and GAB information for support analysis

You must run the `getcomms` script to gather LLT and GAB information when you encounter issues with LLT and GAB. The `getcomms` script also collects core dump and stack traces along with the LLT and GAB information.

To gather LLT and GAB information for support analysis

- 1 If you had changed the default value of the `GAB_FFDC_LOGDIR` parameter, you must again export the same variable before you run the `getcomms` script.

See “GAB message logging” on page 187.

- 2 Run the following command to gather information:

```
# /opt/VRTSgab/getcomms
```

The script uses `remsh` by default. Make sure that you have configured passwordless `remsh`. If you have passwordless `ssh` between the cluster nodes, you can use the `-ssh` option. To gather information on the node that you run the command, use the `-local` option.

Troubleshoot and fix the issue.

If the issue cannot be fixed, then contact Symantec technical support with the file `/tmp/commslog.time_stamp.tar` that the `getcomms` script generates.

Gathering IMF information for support analysis

You must run the `getimf` script to gather information when you encounter issues with IMF (Intelligent Monitoring Framework).

To gather IMF information for support analysis

- ◆ Run the following command on each node:

```
# /opt/VRTSamf/bin/getimf
```

Troubleshoot and fix the issue.

If the issue cannot be fixed, then contact Symantec technical support with the file that the `getimf` script generates.

SF Oracle RAC log files

[Table 3-1](#) lists the various log files and their location. The log files contain useful information for identifying issues and resolving them.

Table 3-1 List of log files

Log file	Location	Description
Oracle installation error log	<code>oraInventory_path\ /logs\ installActionsdate_time.log</code>	<p>Contains errors that occurred during Oracle RAC installation. It clarifies the nature of the error and when it occurred during the installation.</p> <p>Note: Verify if there are any installation errors logged in this file, since they may prove to be critical errors. If there are any installation problems, send this file to Tech Support for debugging the issue.</p>
Oracle alert log	<p>For Oracle RAC 10g:</p> <code>\$ORACLE_HOME/admin/db_name\ bdump/alert_instance_name.log</code> <p>For Oracle RAC 11g:</p> <code>\$ORACLE_BASE/diag/rdbms/db_name\ instance_name/trace/alert_instance_name.log</code> <p>The log path is configurable.</p>	<p>Contains messages and errors reported by database operations.</p>
VCS engine log file	<code>/var/VRTSvcs/log/engine_A.log</code>	<p>Contains all actions performed by the high availability daemon <code>had</code>.</p> <p>Note: Verify if there are any CVM or PrivNIC errors logged in this file, since they may prove to be critical errors.</p>
CVM log files	<code>/var/adm/vx/cmdlog /var/adm/vx/ddl.log /var/adm/vx/translog /var/adm/vx/dmpevents.log /var/VRTSvcs/log/engine_A.log</code>	<p>The <code>cmdlog</code> file contains the list of CVM commands.</p> <p>For more information on collecting important CVM logs: See “Collecting important CVM logs” on page 184.</p>

Table 3-1 List of log files (*continued*)

Log file	Location	Description
VCS agent log files	<p><code>/var/VRTSvcs/log/agenttype_A.log</code> where <i>agenttype</i> is the type of the VCS agent.</p> <p>For example, the log files for the CFS agent can be located at:</p> <p><code>/var/VRTSvcs/log/CFSMount_A.log</code></p>	<p>Contains messages and errors related to the agent functions.</p> <p>For more information, see the <i>Veritas Storage Foundation Administrator's Guide</i>.</p>
OS system log	<code>/var/adm/syslog/syslog.log</code>	Contains messages and errors arising from operating system modules and drivers.
I/O fencing kernel logs	<p><code>/var/VRTSvcs/log/vxfen/vxfen.log</code></p> <p>Obtain the logs by running the following command:</p> <pre># /opt/VRTSvcs/vxfen/bin/\ vxfendebug -p</pre>	Contains messages, errors, or diagnostic information for I/O fencing.
VCSMM log files	<code>/var/VRTSvcs/log/vcsmmconfig.log</code>	Contains messages, errors, or diagnostic information for VCSMM.

Collecting important CVM logs

You need to stop and restart the cluster to collect detailed CVM TIME_JOIN messages.

To collect detailed CVM TIME_JOIN messages

1 On all the nodes in the cluster, perform the following steps.

- Edit the `/opt/VRTSvcs/bin/CVMcluster/online` script.

Insert the '-T' option to the following string.

Original string: `clust_run=`$VXCLUSTADM -m vcs -t $TRANSPORT startnode 2> $CVM_ERR_FILE``


```
Modified string: clust_run=`LANG=C LC_MESSAGES=C $VXCLUSTADM -m
vcs -t $TRANSPORT startnode 2> $CVM_ERR_FILE`
```

2 Stop the cluster.

```
# hastop -all
```

3 Start the cluster.

```
# hastart
```

At this point, CVM TIME_JOIN messages display in the
 /var/adm/syslog/syslog.logfile and on the console.

You can also enable vxconfigd daemon logging as follows:

```
# vxdctl debug 9 /var/adm/vx/vxconfigd_debug.out
```

The debug information that is enabled is accumulated in the system console log and in the text file /var/adm/vx/vxconfigd_debug.out. '9' represents the level of debugging. '1' represents minimal debugging. '9' represents verbose output.

Caution: Turning on vxconfigd debugging degrades VxVM performance. Use vxconfigd debugging with discretion in a production environment.

To disable vxconfigd debugging:

```
# vxdctl debug 0
```

The CVM kernel message dump can be collected on a live node as follows:

```
# /etc/vx/diag.d/kmsgdump -k 2000 > \
/var/adm/vx/kmsgdump.out
```

About SF Oracle RAC kernel and driver messages

SF Oracle RAC drivers such as GAB print messages to the console if the kernel and driver messages are configured to be displayed on the console. Make sure that the kernel and driver messages are logged to the console.

For details on how to configure console messages, see the operating system documentation.

VCS message logging

VCS generates two types of logs: the engine log and the agent log. Log file names are appended by letters. Letter A indicates the first log file, B the second, C the third, and so on.

The engine log is located at `/var/VRTSvcs/log/engine_A.log`. The format of engine log messages is:

Timestamp (Year/MM/DD) | Mnemonic | Severity | UMI | Message Text

- *Timestamp*: the date and time the message was generated.
- *Mnemonic*: the string ID that represents the product (for example, VCS).
- *Severity*: levels include CRITICAL, ERROR, WARNING, NOTICE, and INFO (most to least severe, respectively).
- *UMI*: a unique message ID.
- *Message Text*: the actual message generated by VCS.

A typical engine log resembles:

```
2011/07/10 16:08:09 VCS INFO V-16-1-10077 Received new
cluster membership
```

The agent log is located at `/var/VRTSvcs/log/<agent>.log`. The format of agent log messages resembles:

Timestamp (Year/MM/DD) | Mnemonic | Severity | UMI | Agent Type | Resource Name | Entry Point | Message Text

A typical agent log resembles:

```
2011/07/10 10:38:23 VCS WARNING V-16-2-23331
Oracle:VRT:monitor:Open for ora_lgwr failed, setting
cookie to null.
```

Note that the logs on all nodes may not be identical because

- VCS logs local events on the local nodes.
- All nodes may not be running when an event occurs.

VCS prints the warning and error messages to STDERR.

If the VCS engine, Command Server, or any of the VCS agents encounter some problem, then First Failure Data Capture (FFDC) logs are generated and dumped along with other core dumps and stack traces to the following location:

- For VCS engine: `$VCS_DIAG/diag/had`
- For Command Server: `$VCS_DIAG/diag/CmdServer`

- For VCS agents: `$VCS_DIAG/diag/agents/type`, where *type* represents the specific agent type.

The default value for variable `$VCS_DIAG` is `/var/VRTSvcs/`.

If the debug logging is not turned on, these FFDC logs are useful to analyze the issues that require professional support.

GAB message logging

If GAB encounters some problem, then First Failure Data Capture (FFDC) logs are also generated and dumped.

When you have configured GAB, GAB also starts a GAB logging daemon (`/opt/VRTSgab/gablogd`). GAB logging daemon is enabled by default. You can change the value of the GAB tunable parameter `gab_ibuf_count` to disable the GAB logging daemon.

See [“About GAB load-time or static tunable parameters”](#) on page 238.

This GAB logging daemon collects the GAB related logs when a critical events such as an iofence or failure of the master of any GAB port occur, and stores the data in a compact binary form. You can use the `gabread_ffdc` utility as follows to read the GAB binary log files:

```
/opt/VRTSgab/gabread_ffdc binary_logs_files_location
```

You can change the values of the following environment variables that control the GAB binary log files:

- **GAB_FFDC_MAX_INDX:** Defines the maximum number of GAB binary log files. The GAB logging daemon collects the defined number of log files each of eight MB size. The default value is 20, and the files are named `gablog.1` through `gablog.20`. At any point in time, the most recent file is the `gablog.1` file.
- **GAB_FFDC_LOGDIR:** Defines the log directory location for GAB binary log files

The default location is:

```
/var/adm/gab_ffdc
```

Note that the `gablog` daemon writes its log to the `glgd_A.log` and `glgd_B.log` files in the same directory.

You can either define these variables in the following GAB startup file or use the `export` command. You must restart GAB for the changes to take effect.

```
/etc/rc.config.d/gabconf
```

About debug log tags usage

The following table illustrates the use of debug tags:

Entity	Debug logs used
Agent functions	DBG_1 to DBG_21
Agent framework	DBG_AGTRACE DBG_AGDEBUG DBG_AGINFO
Icmp agent	DBG_HBFW_TRACE DBG_HBFW_DEBUG DBG_HBFW_INFO
HAD	DBG_AGENT (for agent-related debug logs) DBG_ALERTS (for alert debug logs) DBG_CTEAM (for GCO debug logs) DBG_GAB, DBG_GABIO (for GAB debug messages) DBG_GC (for displaying global counter with each log message) DBG_INTERNAL (for internal messages) DBG_IPM (for Inter Process Messaging) DBG_JOIN (for Join logic) DBG_LIC (for licensing-related messages) DBG_NTEVENT (for NT Event logs) DBG_POLICY (for engine policy) DBG_RSM (for RSM debug messages) DBG_TRACE (for trace messages) DBG_SECURITY (for security-related messages) DBG_LOCK (for debugging lock primitives) DBG_THREAD (for debugging thread primitives) DBG_HOSTMON (for HostMonitor debug logs)

Enabling debug logs for agents

This section describes how to enable debug logs for VCS agents.

To enable debug logs for agents

- 1 Set the configuration to read-write:

```
# haconf -makerw
```

- 2 Enable logging and set the desired log levels. The following example depicts the command for the IPMultiNIC resource type.

```
# hatype -modify IPMultiNIC LogDbg DBG_1 DBG_2 DBG_4 DBG_21
```

See the description of the LogDbg attribute for more information.

- 3 For script-based agents, run the `halog` command to add the messages to the engine log:

```
# halog -addtags DBG_1 DBG_2 DBG_4 DBG_21
```

- 4 Save the configuration.

```
# haconf -dump -makero
```

If `DBG_AGDEBUG` is set, the agent framework logs for an instance of the agent appear in the agent log on the node on which the agent is running.

Enabling debug logs for the VCS engine

You can enable debug logs for the VCS engine, VCS agents, and HA commands in two ways:

- To enable debug logs at run-time, use the `halog -addtags` command.
- To enable debug logs at startup, use the `VCS_DEBUG_LOG_TAGS` environment variable. You must set the `VCS_DEBUG_LOG_TAGS` before you start HAD or before you run HA commands.

Examples:

```
# export VCS_DEBUG_LOG_TAGS="DBG_TRACE DBG_POLICY"
# hstart
```

```
# export VCS_DEBUG_LOG_TAGS="DBG_AGINFO DBG_AGDEBUG DBG_AGTRACE"
# hstart
```

```
# export VCS_DEBUG_LOG_TAGS="DBG_IPM"
# hagrps -list
```

Note: Debug log messages are verbose. If you enable debug logs, log files might fill up quickly.

Enabling debug logs for IMF

Run the following commands to enable additional debug logs for Intelligent Monitoring Framework (IMF). The messages get logged in the agent-specific log file `/var/VRTSvcs/log/agentname_A.log`.

To enable additional debug logs

- 1 For Process, Mount, and Application agents:

```
# hatype -modify agentname LogDbg
DBG_AGDEBUG DBG_AGTRACE DBG_AGINFO DBG_1 DBG_2
DBG_3 DBG_4 DBG_5 DBG_6 DBG_7
```

- 2 For Oracle and Netlsnr agents:

```
# hatype -modify agentname LogDbg
DBG_AGDEBUG DBG_AGTRACE DBG_AGINFO DBG_1 DBG_2
DBG_3 DBG_4 DBG_5 DBG_6 DBG_7
DBG_8 DBG_9 DBG_10
```

- 3 For CFSMount agent:

```
# hatype -modify agentname LogDbg
DBG_AGDEBUG DBG_AGTRACE DBG_AGINFO DBG_1 DBG_2
DBG_3 DBG_4 DBG_5 DBG_6 DBG_7
DBG_8 DBG_9 DBG_10 DBG_11 DBG_12
DBG_13 DBG_14 DBG_15 DBG_16
DBG_17 DBG_18 DBG_19 DBG_20 DBG_21
```

- 4 For CVMvxconfigd agent, you do not have to enable any additional debug logs.
- 5 For AMF driver in-memory trace buffer:

```
# amfconfig -S errlevel all all
```

If you had enabled AMF driver in-memory trace buffer, you can view the additional logs using the `amfconfig -p dbglog` command.

Message catalogs

VCS includes multilingual support for message catalogs. These binary message catalogs (BMCs), are stored in the following default locations. The variable *language* represents a two-letter abbreviation.

```
/opt/VRTS/messages/language/module_name
```

The VCS command-line interface displays error and success messages in VCS-supported languages. The `hamsg` command displays the VCS engine logs in VCS-supported languages.

The BMCs are:

<code>gcoconfig.bmc</code>	gcoconfig messages
<code>VRTSvcsHbfw.bmc</code>	Heartbeat framework messages
<code>VRTSvcsTriggers.bmc</code>	VCS trigger messages
<code>VRTSvcsWac.bmc</code>	Wide-area connector process messages
<code>vxfen*.bmc</code>	Fencing messages
<code>gab.bmc</code>	GAB command-line interface messages
<code>hagetcf.bmc</code>	hagetcf messages
<code>llt.bmc</code>	LLT command-line interface messages
<code>VRTSvcsAgfw.bmc</code>	Agent framework messages
<code>VRTSvcsAlerts.bmc</code>	VCS alert messages
<code>VRTSvcsApi.bmc</code>	VCS API messages
<code>VRTSvcsCommon.bmc</code>	Common modules messages
<code>VRTSvcsHad.bmc</code>	VCS engine (HAD) messages
<code>VRTSvcsplatformAgent.bmc</code>	VCS bundled agent messages
<code>VRTSvcsplatformagent_<i>name</i>.bmc</code>	VCS enterprise agent messages

What to do if you see a licensing reminder

In this release, you can install without a license key. In order to comply with the End User License Agreement, you must either install a license key or make the host managed by a Management Server. If you do not comply with these terms within 60 days, the following warning messages result:

```
WARNING V-365-1-1 This host is not entitled to run Veritas Storage
Foundation/Veritas Cluster Server.As set forth in the End User
License Agreement (EULA) you must complete one of the two options
set forth below. To comply with this condition of the EULA and
stop logging of this message, you have <nn> days to either:
- make this host managed by a Management Server (see
  http://go.symantec.com/sfhakeyless for details and free download),
  or
- add a valid license key matching the functionality in use on this host
  using the command 'vxlicinst'
```

To comply with the terms of the EULA, and remove these messages, you must do one of the following within 60 days:

- Install a valid license key corresponding to the functionality in use on the host. After you install the license key, you must validate the license key using the following command:

```
# /opt/VRTS/bin/vxlicrep
```

- Continue with keyless licensing by managing the server or cluster with a management server.

For more information about keyless licensing, see the following URL:
<http://go.symantec.com/sfhakeyless>

Restarting the installer after a failed connection

If an installation is killed because of a failed connection, you can restart the installer to resume the installation. The installer detects the existing installation. The installer prompts you whether you want to resume the installation. If you resume the installation, the installation proceeds from the point where the installation failed.

Installer cannot create UUID for the cluster

The installer displays the following error message if the installer cannot find the `uuidconfig.pl` script before it configures the UUID for the cluster:

```
Couldn't find uuidconfig.pl for uuid configuration,
please create uuid manually before start vcs
```

You may see the error message during SF Oracle RAC configuration, upgrade, or when you add a node to the cluster using the installer.

Workaround: To start SF Oracle RAC, you must run the `uuidconfig.pl` script manually to configure the UUID on each cluster node.

To configure the cluster UUID when you create a cluster manually

- ◆ On one node in the cluster, perform the following command to populate the cluster UUID on each node in the cluster.

```
# /opt/VRTSvcs/bin/uuidconfig.pl -clus -configure nodeA
nodeB ... nodeN
```

Where nodeA, nodeB, through nodeN are the names of the cluster nodes.

Troubleshooting SF Oracle RAC pre-installation check failures

Table 3-2 provides guidelines for resolving failures that may be reported when you start the SF Oracle RAC installation program.

Table 3-2 Troubleshooting pre-installation check failures

Error message	Cause	Resolution
Checking system communication...Failed	Passwordless SSH or remsh communication is not set up between the systems in the cluster.	Set up passwordless SSH or RSH communication between the systems in the cluster. For instructions, see the section "Setting up inter-system communication" in the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> .
Checking release compatibility...Failed	The current system architecture or operating system version is not supported. For example, 32-bit architectures are not supported in this release.	Make sure that the systems meet the hardware and software criteria required for the release. For information, see the chapter "Requirements and planning" in the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> .

Table 3-2 Troubleshooting pre-installation check failures (*continued*)

Error message	Cause	Resolution
Checking installed product	For information on the messages displayed for this check and the appropriate resolutions: See Table 3-3 on page 194.	For information on the messages displayed for this check and the appropriate resolutions: See Table 3-3 on page 194.
Checking platform version	The version of the operating system installed on the system is not supported.	Make sure that the systems meet the software criteria required for the release. For information, see the chapter "Requirements and planning" in the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> .
Performing product prechecks	One of the following checks failed: <ul style="list-style-type: none"> ■ Time synchronization ■ CPU speed checks ■ System architecture checks ■ OS patch level checks 	For information on resolving these issues:

Table 3-3 Checking installed product - messages

Message	Resolution
Entered systems have different products installed: <code>prod_name-prod_ver-sys_name</code> Systems running different products must be upgraded independently.	Two or more systems specified have different products installed. For example, SFCFSHA is installed on some systems while SF Oracle RAC is installed on other systems. You need to upgrade the systems separately. For instructions on upgrading SF Oracle RAC, see the chapter "Upgrading SF Oracle RAC" in the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> . For instructions on upgrading other installed products, see the appropriate product documentation.

Table 3-3 Checking installed product - messages (*continued*)

Message	Resolution
<p>Entered systems have different versions of \$prod_name installed: <i>prod_name-prod_ver-sys_name</i>. Systems running different product versions must be upgraded independently.</p>	<p>Two or more systems specified have different versions of SF Oracle RAC installed. For example, SF Oracle RAC 6.0.1 is installed on some systems while SF Oracle RAC 6.0 is installed on other systems.</p> <p>You need to upgrade the systems separately.</p> <p>For instructions on upgrading SF Oracle RAC, see the chapter "Upgrading SF Oracle RAC" in the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i>.</p>
<p><i>prod_name</i> is installed. Upgrading <i>prod_name</i> directly to <i>another_prod_name</i> is not supported.</p>	<p>A product other than SF Oracle RAC is installed. For example, SFCFSHA 6.0.1 is installed on the systems.</p> <p>The SF Oracle RAC installer does not support direct upgrade from other products to SF Oracle RAC. You need to first upgrade to version 6.0.1 of the installed product, then upgrade to SF Oracle RAC 6.0.1.</p> <p>For instructions on upgrading to SF Oracle RAC, see the section "Upgrading from Storage Foundation products to SF Oracle RAC" in the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i>.</p>
<p><i>prod_name</i> is installed. <i>component_prod_name</i> is part of <i>prod_name</i>. Upgrading only <i>component_prod_name</i> version may partially upgrade the installed depots on the systems.</p>	<p>If a previous version of SF Oracle RAC is installed, the installer supports partial product upgrade. For example, you can upgrade VCS in the stack to version 6.0.1. If you want to upgrade the complete SF Oracle RAC stack later, you can run the <code>installsfrac</code> program.</p>

Troubleshooting LLT health check warning messages

[Table 3-4](#) lists the warning messages displayed during the health check and corresponding recommendations for resolving the issues.

Table 3-4 Troubleshooting LLT warning messages

Warning	Possible causes	Recommendation
Warning: OS timer is not called for <i>num</i> seconds	CPU and memory consumption on the node is high.	Check for applications that may be throttling the CPU and memory resources on the node.
Warning: Kernel failed to allocate memory <i>num</i> time(s)	The available memory on the node is insufficient.	Reduce memory consumption on the node and free up allocated memory.
Flow-control occurred <i>num</i> time(s) and back-enabled <i>num</i> time(s) on port <i>port number</i> for node <i>node number</i>	<ul style="list-style-type: none"> ■ The network bandwidth between the local node and the peer node (<i>node number</i>) is insufficient. ■ The CPU and memory consumption on the peer node (<i>node number</i>) is very high. 	<ul style="list-style-type: none"> ■ Allocate more bandwidth for communication between the local node and the peer node. ■ Check for applications that may be throttling the CPU and memory resources on the node.
Warning: Connectivity with node <i>node id</i> on link <i>link id</i> is flaky <i>num</i> time(s). The distribution is: (0-4 s) <num> (4-8 s) <num> (8-12 s) <num> (12-16 s) <num> (>=16 s)	The private interconnect between the local node and the peer node is unstable.	Check the connectivity between the local node and the peer node. Replace the link, if needed.
One or more link connectivity with peer node(s) <i>node name</i> is in trouble.	The private interconnects between the local node and peer node may not have sufficient bandwidth.	Check the private interconnects between the local node and the peer node.
Link connectivity with <i>node id</i> is on only one link. Symantec recommends configuring a minimum of 2 links.	<ul style="list-style-type: none"> ■ One of the configured private interconnects is non-operational. ■ Only one private interconnect has been configured under LLT. 	<ul style="list-style-type: none"> ■ If the private interconnect is faulty, replace the link. ■ Configure a minimum of two private interconnects.

Table 3-4 Troubleshooting LLT warning messages (*continued*)

Warning	Possible causes	Recommendation
<p>Only one link is configured under LLT. Symantec recommends configuring a minimum of 2 links.</p>	<ul style="list-style-type: none"> ■ One of the configured private interconnects is non-operational. ■ Only one private interconnect has been configured under LLT. 	<ul style="list-style-type: none"> ■ If the private interconnect is faulty, replace the link. ■ Configure a minimum of two private interconnects.
<p>Retransmitted % percentage of total transmitted packets.</p> <p>Sent % percentage of total transmitted packet when no link is up.</p> <p>% percentage of total received packets are with bad checksum.</p> <p>% percentage of total received packets are out of window.</p> <p>% percentage of total received packets are misaligned.</p>	<ul style="list-style-type: none"> ■ The network interface card or the network links are faulty. 	<ul style="list-style-type: none"> ■ Verify the network connectivity between nodes. Check the network interface card for anomalies.
<p>% percentage of total received packets are with DLPI error.</p>	<p>The DLPI driver or NIC may be faulty or corrupted.</p>	<p>Check the DLPI driver and NIC for anomalies.</p>

Table 3-4 Troubleshooting LLT warning messages (*continued*)

Warning	Possible causes	Recommendation
<pre>% per of total transmitted packets are with large xmit latency (>16ms) for port port id %per received packets are with large recv latency (>16ms) for port port id.</pre>	<p>The CPU and memory consumption on the node may be high or the network bandwidth is insufficient.</p>	<ul style="list-style-type: none"> ■ Check for applications that may be throttling CPU and memory usage. ■ Make sure that the network bandwidth is adequate for facilitating low-latency data transmission.
<p>LLT is not running.</p>	<p>SF Oracle RAC is not configured properly.</p>	<p>Reconfigure SF Oracle RAC.</p> <p>For instructions, see the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i>.</p>

Troubleshooting LMX health check warning messages in SF Oracle RAC clusters

[Table 3-5](#) lists the warning messages displayed during the health check and corresponding recommendations for resolving the issues.

Table 3-5 Troubleshooting LMX warning messages

Warning	Possible causes	Recommendation
<pre>LMX is not running. This warning is not applicable for Oracle llg running cluster. VCSMM is not running.</pre>	<p>The possible causes are:</p> <ul style="list-style-type: none"> ■ SF Oracle RAC is not configured properly. ■ SF Oracle RAC is not running. 	<p>Depending on the cause of failure, perform one of the following tasks to resolve the issue:</p> <ul style="list-style-type: none"> ■ Reconfigure SF Oracle RAC. ■ Restart SF Oracle RAC. <p>For instructions: See “Starting or stopping SF Oracle RAC on each node” on page 97.</p>

Table 3-5 Troubleshooting LMX warning messages (*continued*)

Warning	Possible causes	Recommendation
<p>Oracle is not linked to the Symantec LMX library. This warning is not applicable for Oracle 11g running cluster.</p> <p>Oracle is linked to Symantec LMX Library, but LMX is not running.</p> <p>Oracle is not linked to Symantec VCSMM library.</p> <p>Oracle is linked to Symantec VCSMM Library, but VCSMM is not running.</p>	<p>The Oracle RAC libraries are not linked with the SF Oracle RAC libraries.</p>	<p>Relink the Oracle RAC libraries with SF Oracle RAC.</p> <p>For instructions, see the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i>.</p>

Troubleshooting I/O fencing

The following sections discuss troubleshooting the I/O fencing problems. Review the symptoms and recommended solutions.

SCSI reservation errors during bootup

When restarting a node of an SF Oracle RAC cluster, SCSI reservation errors may be observed such as:

```
date system name kernel: scsi3 (0,0,6) : RESERVATION CONFLICT
```

This message is printed for each disk that is a member of any shared disk group which is protected by SCSI-3 PR I/O fencing. This message may be safely ignored.

The vxfcntl utility fails when SCSI TEST UNIT READY command fails

While running the vxfcntl utility, you may see a message that resembles as follows:

```
Issuing SCSI TEST UNIT READY to disk reserved by other node
FAILED.
Contact the storage provider to have the hardware configuration
fixed.
```

The disk array does not support returning success for a SCSI TEST UNIT READY command when another host has the disk reserved using SCSI-3 persistent reservations. This happens with the Hitachi Data Systems 99XX arrays if bit 186 of the system mode option is not enabled.

Node is unable to join cluster while another node is being ejected

A cluster that is currently fencing out (ejecting) a node from the cluster prevents a new node from joining the cluster until the fencing operation is completed. The following are example messages that appear on the console for the new node:

```
...VxFEN ERROR V-11-1-25 ... Unable to join running cluster
since cluster is currently fencing
a node out of the cluster.
```

If you see these messages when the new node is booting, the vxfcntl startup script on the node makes up to five attempts to join the cluster.

To manually join the node to the cluster when I/O fencing attempts fail

- ◆ If the vxfcntl script fails in the attempts to allow the node to join the cluster, restart vxfcntl driver with the command:

```
# /sbin/init.d/vxfen stop
# /sbin/init.d/vxfen start
```

If the command fails, restart the new node.

System panics to prevent potential data corruption

When a node experiences a split-brain condition and is ejected from the cluster, it panics and displays the following console message:

```
VXFEN:vxfcntl_plat_panic: Local cluster node ejected from cluster to
prevent potential data corruption.
```


A node experiences the split-brain condition when it loses the heartbeat with its peer nodes due to failure of all private interconnects or node hang. Review the behavior of I/O fencing under different scenarios and the corrective measures to be taken.

See [“How I/O fencing works in different event scenarios”](#) on page 49.

Cluster ID on the I/O fencing key of coordinator disk does not match the local cluster’s ID

If you accidentally assign coordinator disks of a cluster to another cluster, then the fencing driver displays an error message similar to the following when you start I/O fencing:

```
000068 06:37:33 2bdd5845 0 ... 3066 0 VXFEN WARNING V-11-1-56
Coordinator disk has key with cluster id 48813
which does not match local cluster id 57069
```

The warning implies that the local cluster with the cluster ID 57069 has keys. However, the disk also has keys for cluster with ID 48813 which indicates that nodes from the cluster with cluster id 48813 potentially use the same coordinator disk.

You can run the following commands to verify whether these disks are used by another cluster. Run the following commands on one of the nodes in the local cluster. For example, on `sys1`:

```
sys1> # lltstat -C
57069

sys1> # cat /etc/vxfentab
/dev/vx/rdmp/disk_7
/dev/vx/rdmp/disk_8
/dev/vx/rdmp/disk_9

sys1> # vxfenadm -s /dev/vx/rdmp/disk_7
Reading SCSI Registration Keys...
Device Name: /dev/vx/rdmp/disk_7
Total Number Of Keys: 1
key[0]:
[Numeric Format]: 86,70,48,49,52,66,48,48
[Character Format]: VFBEAD00
[Node Format]: Cluster ID: 48813 Node ID: 0 Node Name: unknown
```

Where `disk_7`, `disk_8`, and `disk_9` represent the disk names in your setup.

Recommended action: You must use a unique set of coordinator disks for each cluster. If the other cluster does not use these coordinator disks, then clear the keys using the `vxfcntlpre` command before you use them as coordinator disks in the local cluster.

See [“About the vxfcntlpre utility”](#) on page 139.

Fencing startup reports preexisting split-brain

The `vxfen` driver functions to prevent an ejected node from rejoining the cluster after the failure of the private network links and before the private network links are repaired.

For example, suppose the cluster of system 1 and system 2 is functioning normally when the private network links are broken. Also suppose system 1 is the ejected system. When system 1 restarts before the private network links are restored, its membership configuration does not show system 2; however, when it attempts to register with the coordinator disks, it discovers system 2 is registered with them. Given this conflicting information about system 2, system 1 does not join the cluster and returns an error from `vxfenconfig` that resembles:

```
vxfenconfig: ERROR: There exists the potential for a preexisting
split-brain. The coordinator disks list no nodes which are in
the current membership. However, they also list nodes which are
not in the current membership.
```

```
I/O Fencing Disabled!
```

Note: During the system boot, because the HP-UX `rc` sequencer redirects the `stderr` of all `rc` scripts to the file `/etc/rc.log`, the error messages will not be printed on the console. It will be logged in the `/etc/rc.log` file.

Also, the following information is displayed on the console:

```
<date> <system name> vxfen: WARNING: Potentially a preexisting
<date> <system name> split-brain.
<date> <system name> Dropping out of cluster.
<date> <system name> Refer to user documentation for steps
<date> <system name> required to clear preexisting split-brain.
<date> <system name>
<date> <system name> I/O Fencing DISABLED!
<date> <system name>
<date> <system name> gab: GAB:20032: Port b closed
```

Note: If `syslogd` is configured with the `-D` option, then the informational message will not be printed on the console. The messages will be logged in the system buffer. The system buffer can be read with the `dmesg` command.

However, the same error can occur when the private network links are working and both systems go down, system 1 restarts, and system 2 fails to come back up. From the view of the cluster from system 1, system 2 may still have the registrations on the coordination points.

Assume the following situations to understand preexisting split-brain in server-based fencing:

- There are three CP servers acting as coordination points. One of the three CP servers then becomes inaccessible. While in this state, one client node leaves the cluster, whose registration cannot be removed from the inaccessible CP server. When the inaccessible CP server restarts, it has a stale registration from the node which left the SF Oracle RAC cluster. In this case, no new nodes can join the cluster. Each node that attempts to join the cluster gets a list of registrations from the CP server. One CP server includes an extra registration (of the node which left earlier). This makes the joiner node conclude that there exists a preexisting split-brain between the joiner node and the node which is represented by the stale registration.
- All the client nodes have crashed simultaneously, due to which fencing keys are not cleared from the CP servers. Consequently, when the nodes restart, the `vxfen` configuration fails reporting preexisting split brain.

These situations are similar to that of preexisting split-brain with coordinator disks, where you can solve the problem running the `vxfenclearpre` command. A similar solution is required in server-based fencing using the `cpsadm` command.

See [“Clearing preexisting split-brain condition”](#) on page 203.

Clearing preexisting split-brain condition

Review the information on how the VxFEN driver checks for preexisting split-brain condition.

See [“Fencing startup reports preexisting split-brain”](#) on page 202.

[Table 3-6](#) describes how to resolve a preexisting split-brain condition depending on the scenario you have encountered:

Table 3-6 Recommended solution to clear pre-existing split-brain condition

Scenario	Solution
Actual potential split-brain condition—system 2 is up and system 1 is ejected	<ol style="list-style-type: none"> 1 Determine if system1 is up or not. 2 If system 1 is up and running, shut it down and repair the private network links to remove the split-brain condition. 3 Restart system 1.
Apparent potential split-brain condition—system 2 is down and system 1 is ejected (Disk-based fencing is configured)	<ol style="list-style-type: none"> 1 Physically verify that system 2 is down. Verify the systems currently registered with the coordination points. Use the following command for coordinator disks: <pre># vxfenadm -s all -f /etc/vxfentab</pre> The output of this command identifies the keys registered with the coordinator disks. 2 Clear the keys on the coordinator disks as well as the data disks in all shared disk groups using the <code>vxfcntlclearpre</code> command. The command removes SCSI-3 registrations and reservations. See “About the vxfcntlclearpre utility” on page 139. 3 Make any necessary repairs to system 2. 4 Restart system 2.

Table 3-6 Recommended solution to clear pre-existing split-brain condition
(continued)

Scenario	Solution
<p>Apparent potential split-brain condition—system 2 is down and system 1 is ejected (Server-based fencing is configured)</p>	<p>1 Physically verify that system 2 is down.</p> <p>Verify the systems currently registered with the coordination points.</p> <p>Use the following command for CP servers:</p> <pre># cpsadm -s cp_server -a list_membership -c cluster_name</pre> <p>where <i>cp_server</i> is the virtual IP address or virtual hostname on which CP server is configured, and <i>cluster_name</i> is the VCS name for the SF Oracle RAC cluster (application cluster).</p> <p>The command lists the systems registered with the CP server.</p> <p>2 Clear the keys on the CP servers using the <code>cpsadm</code> command. The <code>cpsadm</code> command clears a registration on a CP server:</p> <pre># cpsadm -s cp_server -a unreg_node -c cluster_name -n nodeid</pre> <p>where <i>cp_server</i> is the virtual IP address or virtual hostname on which the CP server is listening, <i>cluster_name</i> is the VCS name for the SF Oracle RAC cluster, and <i>nodeid</i> specifies the node id of SF Oracle RAC cluster node. Ensure that fencing is not already running on a node before clearing its registration on the CP server.</p> <p>After removing all stale registrations, the joiner node will be able to join the cluster.</p> <p>3 Make any necessary repairs to system 2.</p> <p>4 Restart system 2.</p>

Registered keys are lost on the coordinator disks

If the coordinator disks lose the keys that are registered, the cluster might panic when a cluster reconfiguration occurs.

To refresh the missing keys

- ◆ Use the `vxfsnwap` utility to replace the coordinator disks with the same disks. The `vxfsnwap` utility registers the missing keys during the disk replacement.

See [“Refreshing lost keys on coordinator disks”](#) on page 152.

Replacing defective disks when the cluster is offline

If the disk becomes defective or inoperable and you want to switch to a new diskgroup in a cluster that is offline, then perform the following procedure.

In a cluster that is online, you can replace the disks using the `vxfsnwap` utility.

See [“About the vxfsnwap utility”](#) on page 142.

Review the following information to replace coordinator disk in the coordinator disk group, or to destroy a coordinator disk group.

Note the following about the procedure:

- When you add a disk, add the disk to the disk group `vxfsncoorddg` and retest the group for support of SCSI-3 persistent reservations.
- You can destroy the coordinator disk group such that no registration keys remain on the disks. The disks can then be used elsewhere.

To replace a disk in the coordinator disk group when the cluster is offline

- 1 Log in as superuser on one of the cluster nodes.
- 2 If VCS is running, shut it down:

```
# hstop -all
```

Make sure that the port `h` is closed on all the nodes. Run the following command to verify that the port `h` is closed:

```
# gabconfig -a
```

- 3 Stop the VCSMM driver on each node:

```
# /sbin/init.d/vcsmm stop
```

- 4 Stop I/O fencing on each node:

```
# /sbin/init.d/vxfen stop
```

This removes any registration keys on the disks.

- 5 Import the coordinator disk group. The file `/etc/vxfendg` includes the name of the disk group (typically, `vxfencoordg`) that contains the coordinator disks, so use the command:

```
# vxdg -tfc import `cat /etc/vxfendg`
```

where:

-t specifies that the disk group is imported only until the node restarts.

-f specifies that the import is to be done forcibly, which is necessary if one or more disks is not accessible.

-C specifies that any import locks are removed.

- 6 To remove disks from the disk group, use the VxVM disk administrator utility, `vxdiskadm`.

You may also destroy the existing coordinator disk group. For example:

- Verify whether the coordinator attribute is set to on.

```
# vxdg list vxfencoordg | grep flags: | grep coordinator
```

- Destroy the coordinator disk group.

```
# vxdg -o coordinator destroy vxfencoordg
```

- 7 Add the new disk to the node and initialize it as a VxVM disk.

Then, add the new disk to the `vxfencoordg` disk group:

- If you destroyed the disk group in step 6, then create the disk group again and add the new disk to it.

See the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide* for detailed instructions.

- If the diskgroup already exists, then add the new disk to it.

```
# vxdg -g vxfencoordg -o coordinator adddisk disk_name
```

- 8 Test the recreated disk group for SCSI-3 persistent reservations compliance.

See “[Testing the coordinator disk group using `vxfsentsthdw -c` option](#)” on page 129.

- 9 After replacing disks in a coordinator disk group, deport the disk group:

```
# vxdg deport `cat /etc/vxfendg`
```

10 On each node, start the I/O fencing driver:

```
# /sbin/init.d/vxfen start
```

11 On each node, start the VCSMM driver:

```
# /sbin/init.d/vcsmm start
```

12 Verify that the I/O fencing module has started and is enabled.

```
# gabconfig -a
```

Make sure that port b and port o memberships exist in the output for all nodes in the cluster.

```
# vxfenadm -d
```

Make sure that I/O fencing mode is not disabled in the output.

13 If necessary, restart VCS on each node:

```
# hastart
```

Troubleshooting I/O fencing health check warning messages

[Table 3-7](#) lists the warning messages displayed during the health check and corresponding recommendations for resolving the issues.

Table 3-7 Troubleshooting I/O fencing warning messages

Warning	Possible causes	Recommendation
VxFEN is not running on all nodes in the cluster.	I/O fencing is not enabled in the cluster.	Configure fencing. For instructions, see the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> .
VxFEN is not configured in SCSI3 mode.	<ul style="list-style-type: none"> ■ If you are using disk-based I/O fencing, it is not configured in the SCSI3 mode. 	Start fencing in SCSI3 mode on all nodes in the cluster.

Table 3-7 Troubleshooting I/O fencing warning messages (*continued*)

Warning	Possible causes	Recommendation
<p>VxFEN is running with only one coordinator disk. Loss of this disk will prevent cluster reconfiguration on loss of a node. Symantec recommends configuring a minimum of 3 coordinator disks.</p> <p>VxFEN is running with even number of coordinator disks. There must be odd number of coordinator disks.</p> <p>Replace the disk <i>disk name</i> using OCDR procedure.</p>	<ul style="list-style-type: none"> ■ I/O fencing is enabled with only one coordinator disk. ■ The number of disks configured for I/O fencing are even. 	<ul style="list-style-type: none"> ■ Configure a minimum of three coordinator disks. ■ Add or remove coordinators disks to meet the odd coordinator disk criterion. Use the <code>vxfsnwap</code> utility to add or remove disks. See “About the vxfsnwap utility” on page 142. ■ If the coordinator disk is faulty, replace the disk using the <code>vxfsnwap</code> utility. See “About the vxfsnwap utility” on page 142.
<p>The coordinator disk (<i>disk name</i>) does not have the required key for the local node.</p>	<p>The key may have been accidentally removed.</p>	<p>Refresh the key on the coordinator disk.</p> <p>For instructions: See “Refreshing lost keys on coordinator disks” on page 152.</p>
<p>SCSI3 write-exclusive reservation is missing on shared disk (<i>disk_name</i>)</p>	<p>The SCSI3 reservation is accidentally removed.</p>	<p>Shut down and restart SF Oracle RAC.</p> <p>When CVM restarts, CVM re-registers the key and places the reservation.</p>

Troubleshooting CP server

All CP server operations and messages are logged in the `/var/VRTScps/log` directory in a detailed and easy to read format. The entries are sorted by date and

time. The logs can be used for troubleshooting purposes or to review for any possible security issue on the system that hosts the CP server.

The following files contain logs and text files that may be useful in understanding and troubleshooting a CP server:

- `/var/VRTScps/log/cpserver_[ABC].log`
- `/var/VRTSvcS/log/vcsauthserver.log` (Security related)
- If the `vxcpserv` process fails on the CP server, then review the following diagnostic files:
 - `/var/VRTScps/diag/FFDC_CPS_pid_vxcpserv.log`
 - `/var/VRTScps/diag/stack_pid_vxcpserv.txt`

Note: If the `vxcpserv` process fails on the CP server, these files are present in addition to a core file. VCS restarts `vxcpserv` process automatically in such situations.

The file `/var/VRTSvcS/log/vxfen/vxfend_[ABC].log` contains logs that may be useful in understanding and troubleshooting fencing-related issues on a SF Oracle RAC cluster (client cluster) node.

See [“Troubleshooting issues related to the CP server service group”](#) on page 210.

See [“Checking the connectivity of CP server”](#) on page 211.

See [“Issues during fencing startup on SF Oracle RAC cluster nodes set up for server-based fencing”](#) on page 211.

See [“Issues during online migration of coordination points”](#) on page 212.

Troubleshooting issues related to the CP server service group

If you cannot bring up the CPSSG service group after the CP server configuration, perform the following steps:

- Verify that the CPSSG service group and its resources are valid and properly configured in the VCS configuration.
- Check the VCS engine log (`/var/VRTSvcS/log/engine_[ABC].log`) to see if any of the CPSSG service group resources are **FAULTED**.
- Review the sample dependency graphs to make sure the required resources are configured correctly.

Checking the connectivity of CP server

You can test the connectivity of CP server using the `cpsadm` command.

You must have set the environment variables `CPS_USERNAME` and `CPS_DOMAINTYPE` to run the `cpsadm` command on the SF Oracle RAC cluster (client cluster) nodes.

To check the connectivity of CP server

- ◆ Run the following command to check whether a CP server is up and running at a process level:

```
# cpsadm -s cp_server -a ping_cps
```

where `cp_server` is the virtual IP address or virtual hostname on which the CP server is listening.

Troubleshooting server-based fencing on the SF Oracle RAC cluster nodes

The file `/var/VRTSvcs/log/vxfen/vxfend_[ABC].log` contains logs files that may be useful in understanding and troubleshooting fencing-related issues on a SF Oracle RAC cluster (application cluster) node.

Issues during fencing startup on SF Oracle RAC cluster nodes set up for server-based fencing

Table 3-8 Fencing startup issues on SF Oracle RAC cluster (client cluster) nodes

Issue	Description and resolution
<code>cpsadm</code> command on the SF Oracle RAC cluster gives connection error	<p>If you receive a connection error message after issuing the <code>cpsadm</code> command on the SF Oracle RAC cluster, perform the following actions:</p> <ul style="list-style-type: none"> ■ Ensure that the CP server is reachable from all the SF Oracle RAC cluster nodes. ■ Check that the SF Oracle RAC cluster nodes use the correct CP server virtual IP or virtual hostname and the correct port number. Check the <code>/etc/vxfenmode</code> file. ■ Ensure that the running CP server is using the same virtual IP/virtual hostname and port number.

Table 3-8 Fencing startup issues on SF Oracle RAC cluster (client cluster) nodes (*continued*)

Issue	Description and resolution
Authorization failure	<p>Authorization failure occurs when the CP server's nodes or users are not added in the CP server configuration. Therefore, fencing on the SF Oracle RAC cluster (client cluster) node is not allowed to access the CP server and register itself on the CP server. Fencing fails to come up if it fails to register with a majority of the coordination points.</p> <p>To resolve this issue, add the CP server node and user in the CP server configuration and restart fencing.</p>
Authentication failure	<p>If you had configured secure communication between the CP server and the SF Oracle RAC cluster (client cluster) nodes, authentication failure can occur due to the following causes:</p> <ul style="list-style-type: none"> ■ Symantec Product Authentication Services (AT) is not properly configured on the CP server and/or the SF Oracle RAC cluster. ■ The CP server and the SF Oracle RAC cluster nodes use different root brokers, and trust is not established between the authentication brokers: See “About secure communication between the SF Oracle RAC cluster and CP server” on page 77.

Issues during online migration of coordination points

During online migration of coordination points using the `vxfenmode` utility, the operation is automatically rolled back if a failure is encountered during validation of coordination points from any of the cluster nodes.

Validation failure of the new set of coordination points can occur in the following circumstances:

- The `/etc/vxfenmode.test` file is not updated on all the SF Oracle RAC cluster nodes, because new coordination points on the node were being picked up from an old `/etc/vxfenmode.test` file. The `/etc/vxfenmode.test` file must be updated with the current details. If the `/etc/vxfenmode.test` file is not present, `vxfenmode` copies configuration for new coordination points from the `/etc/vxfenmode` file.
- The coordination points listed in the `/etc/vxfenmode` file on the different SF Oracle RAC cluster nodes are not the same. If different coordination points are listed in the `/etc/vxfenmode` file on the cluster nodes, then the operation fails due to failure during the coordination point snapshot check.
- There is no network connectivity from one or more SF Oracle RAC cluster nodes to the CP server(s).

- Cluster, nodes, or users for the SF Oracle RAC cluster nodes have not been added on the new CP servers, thereby causing authorization failure.

Vxfen service group activity after issuing the vxfenswap command

The Coordination Point agent reads the details of coordination points from the `vxfenconfig -l` output and starts monitoring the registrations on them.

Thus, during `vxfenswap`, when the `vxfenmode` file is being changed by the user, the Coordination Point agent does not move to FAULTED state but continues monitoring the old set of coordination points.

As long as the changes to `vxfenmode` file are not committed or the new set of coordination points are not reflected in `vxfenconfig -l` output, the Coordination Point agent continues monitoring the old set of coordination points it read from `vxfenconfig -l` output in every monitor cycle.

The status of the Coordination Point agent (either ONLINE or FAULTED) depends upon the accessibility of the coordination points, the registrations on these coordination points, and the fault tolerance value.

When the changes to `vxfenmode` file are committed and reflected in the `vxfenconfig -l` output, then the Coordination Point agent reads the new set of coordination points and proceeds to monitor them in its new monitor cycle.

Troubleshooting Cluster Volume Manager in SF Oracle RAC clusters

This section discusses troubleshooting CVM problems.

Restoring communication between host and disks after cable disconnection

If a fiber cable is inadvertently disconnected between the host and a disk, you can restore communication between the host and the disk without restarting.

To restore lost cable communication between host and disk

- 1 Reconnect the cable.
- 2 On all nodes, use the `ioscan -funC disk` command to scan for new disks.

It may take a few minutes before the host is capable of seeing the disk.

- 3 On all nodes, issue the following command to rescan the disks:

```
# vxdisk scandisks
```

- 4 On the master node, reattach the disks to the disk group they were in and retain the same media name:

```
# vxreattach
```

This may take some time. For more details, see `vxreattach(1M)` manual page.

Shared disk group cannot be imported in SF Oracle RAC cluster

If you see a message resembling:

```
vxvm:vxconfigd:ERROR:vold_pgr_register(/dev/vx/rdmp/disk_name):  
local_node_id<0  
Please make sure that CVM and vxfen are configured  
and operating correctly
```

First, make sure that CVM is running. You can see the CVM nodes in the cluster by running the `vxclustadm nidmap` command.

```
# vxclustadm nidmap  
Name          CVM Nid    CM Nid     State  
sys1          1          0          Joined: Master  
sys2          0          1          Joined: Slave
```

This above output shows that CVM is healthy, with system `sys1` as the CVM master. If CVM is functioning correctly, then the output above is displayed when CVM cannot retrieve the node ID of the local system from the `vxfen` driver. This usually happens when port `b` is not configured.

To verify vxfen driver is configured

- ◆ Check the GAB ports with the command:

```
# gabconfig -a
```

Port `b` must exist on the local system.

Error importing shared disk groups in SF Oracle RAC cluster

The following message may appear when importing shared disk group:

```
VxVM vxdg ERROR V-5-1-587 Disk group disk_group_name: import  
failed: No valid disk found containing disk group
```

You may need to remove keys written to the disk.

For information about removing keys written to the disk:

See [“Removing preexisting keys”](#) on page 140.

Unable to start CVM in SF Oracle RAC cluster

If you cannot start CVM, check the consistency between the `/etc/llthosts` and `main.cf` files for node IDs.

You may need to remove keys written to the disk.

For information about removing keys written to the disk:

See [“Removing preexisting keys”](#) on page 140.

CVM group is not online after adding a node to the SF Oracle RAC cluster

The possible causes for the CVM group being offline after adding a node to the cluster are as follows:

- The `cssd` resource is configured as a critical resource in the `cvm` group.
- Other resources configured in the `cvm` group as critical resources are not online.

To resolve the issue if `cssd` is configured as a critical resource

- 1 Log onto one of the nodes in the existing cluster as the root user.
- 2 Configure the `cssd` resource as a non-critical resource in the `cvm` group:

```
# haconf -makerw
# hares -modify cssd Critical 0
# haconf -dump -makero
```

To resolve the issue if other resources in the group are not online

- 1 Log onto one of the nodes in the existing cluster as the root user.
- 2 Bring the resource online:

```
# hares -online resource_name -sys system_name
```

- 3 Verify the status of the resource:

```
# hastatus -resource resource_name
```

- 4 If the resource is not online, configure the resource as a non-critical resource:

```
# haconf -makerw
# hares -modify resource_name Critical 0
# haconf -dump -makero
```

CVMVolDg not online even though CVMCluster is online in SF Oracle RAC cluster

When the CVMCluster resource goes online, then all shared disk groups that have the auto-import flag set are automatically imported. If the disk group import fails for some reason, the CVMVolDg resources fault. Clearing and taking the CVMVolDg type resources offline does not resolve the problem.

To resolve the resource issue

- 1 Fix the problem causing the import of the shared disk group to fail.
- 2 Offline the cvm group containing the resource of type CVMVolDg as well as the service group containing the CVMCluster resource type.
- 3 Bring the cvm group containing the CVMCluster resource online.
- 4 Bring the cvm group containing the CVMVolDg resource online.

Troubleshooting VCSIPC

This section discusses troubleshooting VCSIPC problems.

VCSIPC wait warning messages in Oracle trace/log files

When Gigabit Ethernet interconnections are used, a high load can cause LMX/LLT to flow-control VCSIPC, resulting in warning messages to be reported in the Oracle trace file. The default location for the trace file is \$ORACLE_HOME/rdbms/log; it may have changed if the parameters `background_dump_dest` or `user_dump_dest` have been changed. The messages resemble:

```
VCSIPC wait: WARNING: excessive poll done, 1001 times
VCSIPC wait: WARNING: excessive poll done, 1001 times
```


As a workaround, you can change the LLT lowwater mark, highwater mark, and window values for flow control. Please contact Veritas support for more information about changing these values.

VCSIPC errors in Oracle trace/log files

If you see any VCSIPC errors in the Oracle trace/log files, check the `/var/adm/syslog/syslog.log` file for any LMX error messages.

If you see messages that contain any of the following:

```
. . . out of buffers
. . . out of ports
. . . no minors available
```

See “[About LMX tunable parameters](#)” on page 253.

If you see any VCSIPC warning messages in Oracle trace/log files that resemble:

```
connection invalid
```

or,

```
Reporting communication error with node
```

Check whether the Oracle Real Application Cluster instance on the other system is still running or has been restarted. The warning message indicates that the VCSIPC/LMX connection is no longer valid.

Troubleshooting Oracle

This section discusses troubleshooting Oracle.

Oracle log files

The following Oracle log files are helpful for resolving issues with Oracle components:

- [Table 3-9](#)
- [Table 3-10](#)
- [Table 3-11](#)

Table 3-9 Oracle Clusterware log files

Component	Log file location
Cluster Ready Services Daemon (crsd) Log Files	\$CRS_HOME/log/hostname/crsd
Cluster Synchronization Services (CSS)	\$CRS_HOME/log/hostname/cssd
Event Manager (EVM) information generated by evmd	\$CRS_HOME/log/hostname/evmd
Oracle RAC RACG	\$CRS_HOME/log/hostname/racg \$ORACLE_HOME/log/hostname/racg

Table 3-10 Oracle RAC 11g Release 2 log files

Component	Log file location
Clusterware alert log	\$GRID_HOME/log/<host>/alert<host>.log
Disk Monitor daemon	\$GRID_HOME/log/<host>/diskmon
OCRDUMP, OCRCHECK, OCRCONFIG, CRSCTL	\$GRID_HOME/log/<host>/client
Cluster Time Synchronization Service	\$GRID_HOME/log/<host>/ctssd
Grid Interprocess Communication daemon	\$GRID_HOME/log/<host>/gipcd
Oracle High Availability Services daemon	\$GRID_HOME/log/<host>/ohasd
Cluster Ready Services daemon	\$GRID_HOME/log/<host>/crsd

Table 3-10 Oracle RAC 11g Release 2 log files (*continued*)

Component	Log file location
Grid Plug and Play daemon	\$GRID_HOME/log/<host>/gpnpd
Multicast Domain Name Service daemon	\$GRID_HOME/log/<host>/mdnsd
Event Manager daemon	\$GRID_HOME/log/<host>/evmd
RAC RACG (only used if pre-11.1 database is installed)	\$GRID_HOME/log/<host>/racg
Cluster Synchronization Service daemon	\$GRID_HOME/log/<host>/cssd
Server Manager	\$GRID_HOME/log/<host>/srvm
HA Service Daemon Agent	\$GRID_HOME/log/<host>/agent/ohasd/\ oraagent_oracle11
HA Service Daemon CSS Agent	\$GRID_HOME/log/<host>/agent/ohasd/\ oracssdagent_root
HA Service Daemon ocssd Monitor Agent	\$GRID_HOME/log/<host>/agent/ohasd/\ oracssdmonitor_root
HA Service Daemon Oracle Root Agent	\$GRID_HOME/log/<host>/agent/\ ohasd/orarootagent_root
CRS Daemon Oracle Agent	\$GRID_HOME/log/<host>/agent/\ crsd/oraagent_oracle11
CRS Daemon Oracle Root Agent	\$GRID_HOME/log/<host> agent/\ crsd/orarootagent_root
Grid Naming Service daemon	\$GRID_HOME/log/<host>/gnsd

Table 3-11 Oracle database log file

Component	Log file location
Oracle database	\$ORACLE_BASE/diag/rdbms/database_name/SID/

Oracle Notes

Review the following Oracle notes, when dealing with the following specific Oracle issues:

- 259301.1 Oracle RAC 10g
- 280589.1 Oracle Clusterware installation does not succeed if one or more cluster nodes present are not to be configured for Oracle Clusterware.
- 265769.1 Oracle RAC 10g: Troubleshooting Oracle Clusterware reboots
- 1050693.1 Oracle RAC 11g: Troubleshooting Oracle Clusterware node evictions (reboots)
- 279793.1 How to restore a lost vote disk in Oracle RAC 10g
- 239998.1 Oracle RAC 10g: How to clean up after a failed Oracle Clusterware installation
 - Two items missing in this Oracle note are:
 - Remove the `/var/opt/oracle/ocr.loc` file.
 The file contains the location for the Cluster registry. If this file is not removed then during the next installation the installer will not query for the OCR location and will pick it from this file.
 - If there was a previous 9i Oracle installation, then remove the following file: `/var/opt/oracle/srvConfig.loc`. If this file is present the installer will pick up the Vote disk location from this file and may create the error "the Vote disk should be placed on a shared file system" even before specifying the Vote disk location.
- 1377349.1 How to deconfigure/reconfigure (rebuild OCR) or deinstall Grid Infrastructure
- 330358.1 Oracle Clusterware 10g Release 2/11g Release 1/11g Release 2 Diagnostic Collection Guide

OUI fails to display the nodes names on the cluster configuration screen during the installation of Oracle Clusterware

The Oracle Universal Installer (OUI) fails to display the names of the nodes in the cluster on the "Specify Cluster configuration" screen at the time of installing Oracle Clusterware.

The OUI runs the `/tmp/OraInstall*/*/lsnodes` utility to determine the names of the nodes in the cluster. If the utility fails to retrieve the node information, the node information fails to display in the OUI.

Table 3-12 lists the possible causes of failure and the corresponding resolution.

Table 3-12 OUI failures - Causes and resolution

Cause	Resolution
The <code>/etc/llthosts</code> file is not readable by the Oracle user due to permission issues.	Set the read permission to other users for the <code>/etc/llthosts</code> file.
The Oracle user can not read the following file due to permission issues. HP-UX (IA): <code>/opt/nmapi/nmapi2/lib/hpux64/libnmapi2.so</code>	Set the read permission to other users for the file.

Relinking of VCSMM library fails after upgrading from SF Oracle RAC version 4.1 MP2

After you upgrade from SF Oracle RAC 4.1 MP2, the relinking process may fail with the following message:

```
$CRS_HOME/lib/libskgxn2.so is not a VCSMM library
on one or more node(s) of your cluster.
It should be a symbolic link to
/opt/nmapi/nmapi2/lib/hpux64/libnmapi2.so file,
which must be a VCSMM library file.
```

The process fails because the Veritas `skgxn` library is copied directly to the Oracle Clusterware home directory (`$CRS_HOME/lib`) instead of linking the library in the Oracle Clusterware home directory to the library `/opt/nmapi/nmapi2/lib/hpux64/libnmapi2.so`.

To resolve the issue, create a symbolic link for the library from the Oracle Clusterware home directory to the library `/opt/nmapi/nmapi2/lib/hpux64/libnmapi2.so` as follows:

```
# ln -s /opt/nmapi/nmapi2/lib/hpux64/libnmapi2.so \  
$CRS_HOME/lib/libskgxn2.so
```

After relinking the VCSMM library, relink the ODM library as described in the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide*.

Error when starting an Oracle instance in SF Oracle RAC

If the VCSMM driver (the membership module) is not configured, an error displays while starting the Oracle instance that resembles:

```
ORA-29702: error occurred in Cluster Group Operation
```

To start the VCSMM driver:

```
# /sbin/init.d/vcsmm start
```

Clearing Oracle group faults

If the Oracle group faults, you can clear the faults and bring the group online by running the following commands:

```
# hagrps -clear oracle_grp
```

```
# hagrps -online oracle_grp -sys node_name
```

Oracle log files show shutdown called even when not shutdown manually

The Oracle enterprise agent calls shutdown if monitoring of the Oracle resources fails. On all cluster nodes, review the following VCS and Oracle agent log files for any errors or status:

```
/var/VRTSvcs/log/engine_A.log
```

```
/var/VRTSvcs/log/Oracle_A.log
```

Resolving ASYNCH_IO errors in an SF Oracle RAC cluster

If ASYNCH_IO errors occur during select and update queries on the Oracle database, the workaround involves setting the MLOCK privilege for the dba user.

To set MLOCK privilege for DBA user

- 1 Give the RTPRIO MLOCK RTSCHED privilege to the dba group:

```
# setprivgrp dba RTPRIO MLOCK RTSCHED
```

- 2 Create the `/etc/privgroup` file and add the line:

```
dba RTPRIO MLOCK RTSCHED
```

- 3 Verify the availability of the privilege for the dba group:

```
# /usr/bin/getprivgrp dba
```

Oracle's clusterware processes fail to start

Verify that the Oracle RAC configuration meets the following configuration requirements:

- The correct private IP address is configured on the private link using the PrivNIC or MultiPrivNIC agent.
- The OCR and voting disks shared disk groups are accessible.
- The file systems containing OCR and voting disks are mounted.

Check the CSSD log files to learn more.

For Oracle RAC 10g Release 2:

You can find the CSSD log files at `$CRS_HOME/log/node_name/cssd/*`

For Oracle RAC 11g Release 2:

You can find the CSSD log files at `$GRID_HOME/log/node_name/cssd/*`

Consult the Oracle RAC documentation for more information.

Oracle Clusterware fails after restart

If the Oracle Clusterware fails to start after boot up, check for the occurrence of the following strings in the `/var/adm/syslog/syslog.log` messages file.

String value in the file:

```
Oracle CSSD failure.  
Rebooting for cluster  
integrity
```

Oracle Clusterware may fail due to Oracle CSSD failure. The Oracle CSSD failure may be caused by one of the following events:

- Communication failure occurred and Oracle Clusterware fenced out the node.
- OCR and Vote disk became unavailable.
- `ocssd` was killed manually.
- Killing the `init.cssd` script.

String value in the file:

```
Waiting for file  
system containing
```

The Oracle Clusterware installation is on a shared disk and the `init` script is waiting for that file system to be made available.

String value in the file:

```
Oracle Cluster Ready  
Services disabled by  
corrupt install
```

The following file is not available or has corrupt entries:

```
/var/opt/oracle/scls_scr/  
hostname/root/crsstart
```

String value in the file:

```
OCR initialization  
failed accessing OCR  
device
```

The shared file system containing the OCR is not available and Oracle Clusterware is waiting for it to become available.

Troubleshooting the Virtual IP (VIP) configuration in an SF Oracle RAC cluster

When troubleshooting issues with the VIP configuration, use the following commands and files:

- Check for network problems on all nodes:

```
# /usr/sbin/ifconfig nic_name
```

- Verify the `/etc/hosts` file on each node. Or, make sure that the virtual host name is registered with the DNS server as follows:
- Verify the virtual host name on each node.

```
# ping virtual_host_name
```

- Check the output of the following command:

For Oracle RAC 10g Release 2/Oracle RAC 11g Release 1:

```
$CRS_HOME/bin/crs_stat -t
```

For Oracle RAC 11g Release 2:

```
$GRID_HOME/bin/crsctl stat res -t
```

- On the problem node, use the command:
- For Oracle RAC 10g Release 2/Oracle RAC 11g Release 1:

```
$ $CRS_HOME/bin/srvctl start nodeapps -n node_name
```

For Oracle RAC 11g Release 2:

```
$ $GRID_HOME/bin/srvctl start nodeapps -n node_name
```

Troubleshooting Oracle Clusterware health check warning messages in SF Oracle RAC clusters

Table 3-13 lists the warning messages displayed during the health check and corresponding recommendations for resolving the issues.

Table 3-13 Troubleshooting Oracle Clusterware warning messages

Warning	Possible causes	Recommendation
Oracle Clusterware is not running.	<p>Oracle Clusterware is not started.</p> <p>Oracle Clusterware is waiting for dependencies such as OCR or voting disk or private IP addresses to be available.</p>	<ul style="list-style-type: none"> ■ Bring the cvm group online to start Oracle Clusterware: <pre># hagrpx -online cvm \ -sys system_name</pre> ■ Check the VCS log file <code>/var/VRTSvcs/log/engine_A.log</code> to determine the cause and fix the issue.
No CSSD resource is configured under VCS.	The CSSD resource is not configured under VCS.	<p>Configure the CSSD resource under VCS and bring the resource online.</p> <p>See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i>.</p>

Table 3-13 Troubleshooting Oracle Clusterware warning messages (*continued*)

Warning	Possible causes	Recommendation
The CSSD resource <i>name</i> is not running.	<ul style="list-style-type: none"> ■ VCS is not running. ■ The dependent resources, such as the CFSMount or CVMVoIdg resource for OCR and voting disk are not online. 	<ul style="list-style-type: none"> ■ Start VCS: # hastart ■ Check the VCS log file <code>/var/VRTSvcs/log/engine_A.log</code> to determine the cause and fix the issue.
Mismatch between LLT links <i>llt nics</i> and Oracle Clusterware links <i>crs nics</i> .	The private interconnects used by Oracle Clusterware are not configured over LLT interfaces.	The private interconnects for Oracle Clusterware must use LLT links. Configure the private IP addresses on one of the LLT links. See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> .
Mismatch between Oracle Clusterware links <i>crs nics</i> and PrivNIC links <i>private nics</i> .	The private IP addresses used by Oracle Clusterware are not configured under PrivNIC.	Configure the PrivNIC resource to monitor the private IP address used by Oracle Clusterware. See the <i>Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide</i> .
Mismatch between CRS nodes <i>crs nodes</i> and LLT nodes <i>llt nodes</i> .	The host names configured during the Oracle Clusterware installation are not the same as the host names configured under LLT.	Make sure that the host names configured during the Oracle Clusterware installation are the same as the host names configured under LLT.

Troubleshooting ODM in SF Oracle RAC clusters

This section discusses troubleshooting ODM.

File System configured incorrectly for ODM shuts down Oracle

Linking Oracle RAC with the Veritas ODM libraries provides the best file system performance.

Review the instructions on creating the link and confirming that Oracle uses the libraries. Shared file systems in RAC clusters without ODM libraries linked to Oracle RAC may exhibit slow performance and are not supported.

If ODM cannot find the resources it needs to provide support for cluster file systems, it does not allow Oracle to identify cluster files and causes Oracle to fail at startup.

To verify cluster status, run the following command and review the output:

```
# cat /dev/odm/cluster

cluster status: enabled
```

If the status is "enabled," ODM is supporting cluster files. Any other cluster status indicates that ODM is not supporting cluster files. Other possible values include:

pending	ODM cannot yet communicate with its peers, but anticipates being able to eventually.
failed	ODM cluster support has failed to initialize properly. Check console logs.
disabled	<p>ODM is not supporting cluster files. If you think ODM should be supporting the cluster files:</p> <ul style="list-style-type: none"> ■ Check <code>/dev/odm</code> mount options using <code>mount</code>. If the <code>nocluster</code> option is being used, it can force the <code>disabled</code> cluster support state. ■ Make sure that the <code>VRTSgms</code> (group messaging service) depot is installed. Run the following command: <ul style="list-style-type: none"> # <code>swlist VRTSgms</code> ■ Verify that the <code>gms</code> module is loaded: <ul style="list-style-type: none"> # <code>kcmodule -v vxgms</code> ■ Restart ODM: <ul style="list-style-type: none"> # <code>/sbin/init.d/odm stop</code> # <code>/sbin/init.d/odm start</code>

Prevention and recovery strategies

This chapter includes the following topics:

- [Verification of GAB ports in SF Oracle RAC cluster](#)
- [Examining GAB seed membership](#)
- [Manual GAB membership seeding](#)
- [Evaluating VCS I/O fencing ports](#)
- [Verifying normal functioning of VCS I/O fencing](#)
- [Managing SCSI-3 PR keys in SF Oracle RAC cluster](#)
- [Identifying a faulty coordinator LUN](#)

Verification of GAB ports in SF Oracle RAC cluster

The following ports need to be up on all the nodes of SF Oracle RAC cluster:

- GAB
- I/O fencing
- ODM
- CFS
- VCS ('had')
- vcsmm (membership module for SF Oracle RAC)
- CVM (kernel messaging)

- CVM (vxconfigd)
- CVM (to ship commands from slave node to master node)
- CVM (I/O shipping)

The following command can be used to verify the state of GAB ports:

```
# gabconfig -a
```

GAB Port Memberships

```
Port a gen 7e6e7e05 membership 01
Port b gen 58039502 membership 01
Port d gen 588a7d02 membership 01
Port f gen 1ea84702 membership 01
Port h gen cf430b02 membership 01
Port o gen de8f0202 membership 01
Port u gen de4f0203 membership 01
Port v gen db411702 membership 01
Port w gen cf430b02 membership 01
Port y gen 73f449 membership 01
```

The data indicates that all the GAB ports are up on the cluster having nodes 0 and 1.

For more information on the GAB ports in SF Oracle RAC cluster, see the *Veritas Storage Foundation for Oracle RAC Installation and Configuration Guide*.

Examining GAB seed membership

The number of systems that participate in the cluster is specified as an argument to the `gabconfig` command in `/etc/gabtab`. In the following example, two nodes are expected to form a cluster:

```
# cat /etc/gabtab
/sbin/gabconfig -c -n2
```

GAB waits until the specified number of nodes becomes available to automatically create the port “a” membership. Port “a” indicates GAB membership for an SF Oracle RAC cluster node. Every GAB reconfiguration, such as a node joining or leaving increments or decrements this seed membership in every cluster member node.

A sample port ‘a’ membership as seen in `gabconfig -a` is shown:

```
Port a gen 7e6e7e01 membership 01
```

In this case, 7e6e7e01 indicates the “membership generation number” and 01 corresponds to the cluster’s “node map”. All nodes present in the node map reflects the same membership ID as seen by the following command:

```
# gabconfig -a | grep "Port a"
```

The semi-colon is used as a placeholder for a node that has left the cluster. In the following example, node 0 has left the cluster:

```
# gabconfig -a | grep "Port a"
```

```
Port a gen 7e6e7e04 membership ;1
```

When the last node exits the port “a” membership, there are no other nodes to increment the membership ID. Thus the port “a” membership ceases to exist on any node in the cluster.

When the last and the final system is brought back up from a complete cluster cold shutdown state, the cluster will seed automatically and form port “a” membership on all systems. Systems can then be brought down and restarted in any combination so long as at least one node remains active at any given time.

The fact that all nodes share the same membership ID and node map certifies that all nodes in the node map participates in the same port “a” membership. This consistency check is used to detect “split-brain” and “pre-existing split-brain” scenarios.

Split-brain occurs when a running cluster is segregated into two or more partitions that have no knowledge of the other partitions. The pre-existing network partition is detected when the “cold” nodes (not previously participating in cluster) start and are allowed to form a membership that might not include all nodes (multiple sub-clusters), thus resulting in a potential split-brain.

Note: Symantec I/O fencing prevents data corruption resulting from any split-brain scenarios.

Manual GAB membership seeding

It is possible that one of the nodes does not come up when all the nodes in the cluster are restarted, due to the “minimum seed requirement” safety that is enforced by GAB. Human intervention is needed to safely determine that the other node is in fact not participating in its own mini-cluster.

The following should be carefully validated before manual seeding, to prevent introducing split-brain and subsequent data corruption:

- Verify that none of the other nodes in the cluster have a port “a” membership

- Verify that none of the other nodes have any shared disk groups imported
- Determine why any node that is still running does not have a port “a” membership

Run the following command to manually seed GAB membership:

```
# gabconfig -cx
```

Refer to `gabconfig (1M)` for more details.

Evaluating VCS I/O fencing ports

I/O Fencing (VxFEN) uses a dedicated port that GAB provides for communication across nodes in the cluster. You can see this port as port ‘b’ when `gabconfig -a` runs on any node in the cluster. The entry corresponding to port ‘b’ in this membership indicates the existing members in the cluster as viewed by I/O Fencing.

GAB uses port “a” for maintaining the cluster membership and must be active for I/O Fencing to start.

To check whether fencing is enabled in a cluster, the ‘-d’ option can be used with `vxfsadm (1M)` to display the I/O Fencing mode on each cluster node. Port “b” membership should be present in the output of `gabconfig -a` and the output should list all the nodes in the cluster.

If the GAB ports that are needed for I/O fencing are not up, that is, if port “a” is not visible in the output of `gabconfig -a` command, LLT and GAB must be started on the node.

The following commands can be used to start LLT and GAB respectively:

To start LLT on each node:

```
# /sbin/init.d/llt start
```

If LLT is configured correctly on each node, the console output displays:

```
LLT INFO V-14-1-10009 LLT Protocol available
```

To start GAB, on each node:

```
# /sbin/init.d/gab start
```

If GAB is configured correctly on each node, the console output displays:

```
GAB INFO V-15-1-20021 GAB available
```

```
GAB INFO V-15-1-20026 Port a registration waiting for seed port membership
```


Verifying normal functioning of VCS I/O fencing

It is mandatory to have VCS I/O fencing enabled in SF Oracle RAC cluster to protect against split-brain scenarios. VCS I/O fencing can be assumed to be running normally in the following cases:

- Fencing port 'b' enabled on the nodes

```
# gabconfig -a
```

To verify that fencing is enabled on the nodes:

```
# vxfenadm -d
```

- Registered keys present on the coordinator disks

```
# vxfenadm -s all -f /etc/vxfentab
```

Managing SCSI-3 PR keys in SF Oracle RAC cluster

I/O Fencing places the SCSI-3 PR keys on coordinator LUNs. The format of the key follows the naming convention wherein ASCII "A" is prefixed to the LLT ID of the system that is followed by 7 dash characters.

For example:

node 0 uses A-----

node 1 uses B-----

In an SF Oracle RAC/SF CFS/SF HA environment, VxVM/CVM registers the keys on data disks, the format of which is ASCII "A" prefixed to the LLT ID of the system followed by the characters "PGRxxxx" where 'xxxx' = i such that the disk group is the ith shared group to be imported.

For example: node 0 uses APGR0001 (for the first imported shared group).

In addition to the registration keys, VCS/CVM also installs a reservation key on the data LUN. There is one reservation key per cluster as only one node can reserve the LUN.

See ["About SCSI-3 Persistent Reservations"](#) on page 46.

The following command lists the keys on a data disk group:

```
# vxdg list |grep data
```

```
sys1_data1 enabled,shared,cds 1201715530.28.pushover
```

Select the data disk belonging to sys1_data1:

```
# vxdisk -o alldgs list |grep sys1_data1

clt2d0s2 auto:cdsdisk clt2d0s2 sys1_data1 online shared
clt2d1s2 auto:cdsdisk clt2d1s2 sys1_data1 online shared
clt2d2s2 auto:cdsdisk clt2d2s2 sys1_data1 online shared
```

The following command lists the PR keys:

```
# vxdisk -o listreserve list clt2d0s2
```

```
.....
.....
```

Alternatively, the PR keys can be listed using `vxfenadm` command:

```
# echo "/dev/vx/dmp/clt2d0s2" > /tmp/disk71

# vxfenadm -s all -f /tmp/disk71

Device Name: /dev/vx/dmp/clt2d0s2
Total Number Of Keys: 2
key[0]:
    Key Value [Numeric Format]: 66,80,71,82,48,48,48,52
    Key Value [Character Format]: BPGR0004
key[1]:
    Key Value [Numeric Format]: 65,80,71,82,48,48,48,52
    Key Value [Character Format]: APGR0004
```

Evaluating the number of SCSI-3 PR keys on a coordinator LUN, if there are multiple paths to the LUN from the hosts

The utility `vxfenadm` (1M) can be used to display the keys on the coordinator LUN. The key value identifies the node that corresponds to each key. Each node installs a registration key on all the available paths to the LUN. Thus, the total number of registration keys is the sum of the keys that are installed by each node in the above manner.

See [“About the vxfenadm utility”](#) on page 134.

Detecting accidental SCSI-3 PR key removal from coordinator LUNs

The keys currently installed on the coordinator disks can be read using the following command:

```
# vxfenadm -s all -f /etc/vxfentab
```

There should be a key for each node in the operating cluster on each of the coordinator disks for normal cluster operation. There will be two keys for every node if you have a two-path DMP configuration.

Identifying a faulty coordinator LUN

The utility `vxfcntlshdw (1M)` provided with I/O fencing can be used to identify faulty coordinator LUNs. This utility must be run from any node in the cluster. The coordinator LUN, which needs to be checked, should be supplied to the utility.

See [“About the vxfcntlshdw utility”](#) on page 127.

Tunable parameters

This chapter includes the following topics:

- [About SF Oracle RAC tunable parameters](#)
- [About GAB tunable parameters](#)
- [About LLT tunable parameters](#)
- [About LMX tunable parameters](#)
- [About VXFEN tunable parameters](#)
- [Tuning guidelines for campus clusters](#)

About SF Oracle RAC tunable parameters

Tunable parameters can be configured to enhance the performance of specific SF Oracle RAC features. This chapter discusses how to configure the following SF Oracle RAC tunables:

- GAB
- LLT
- LMX
- VXFEN

Symantec recommends that you do not change the tunable kernel parameters without assistance from Symantec support personnel. Several of the tunable parameters preallocate memory for critical data structures, and a change in their values could increase memory use or degrade performance.

Warning: Do not adjust the SF Oracle RAC tunable parameters for LMX and VXFEN as described below to enhance performance without assistance from Symantec support personnel.

About GAB tunable parameters

GAB provides various configuration and tunable parameters to modify and control the behavior of the GAB module.

These tunable parameters not only provide control of the configurations like maximum possible number of nodes in the cluster, but also provide control on how GAB behaves when it encounters a fault or a failure condition. Some of these tunable parameters are needed when the GAB module is loaded into the system. Any changes to these load-time tunable parameters require either unload followed by reload of GAB module or system reboot. Other tunable parameters (run-time) can be changed while GAB module is loaded, configured, and cluster is running. Any changes to such a tunable parameter will have immediate effect on the tunable parameter values and GAB behavior.

See [“About GAB load-time or static tunable parameters”](#) on page 238.

See [“About GAB run-time or dynamic tunable parameters”](#) on page 240.

About GAB load-time or static tunable parameters

Table 5-1 lists the static tunable parameters in GAB that are used during module load time. Use the `gabconfig -e` command to list all such GAB tunable parameters.

You can modify these tunable parameters only by adding new values in the GAB configuration file. The changes take effect only on reboot or on reload of the GAB module.

Table 5-1 GAB static tunable parameters

GAB parameter	Description	Values (default and range)
gab_numnids	Maximum number of nodes in the cluster	Default: 64 Range: 1-64
gab_numports	Maximum number of ports in the cluster	Default: 32 Range: 1-32

Table 5-1 GAB static tunable parameters (*continued*)

GAB parameter	Description	Values (default and range)
gab_flowctrl	<p>Number of pending messages in GAB queues (send or receive) before GAB hits flow control.</p> <p>This can be overwritten while cluster is up and running with the <code>gabconfig -Q</code> option. Use the <code>gabconfig</code> command to control value of this tunable.</p>	<p>Default: 128</p> <p>Range: 1-1024</p>
gab_logbufsize	GAB internal log buffer size in bytes	<p>Default: 48100</p> <p>Range: 8100-65400</p>
gab_msglogsize	Maximum messages in internal message log	<p>Default: 256</p> <p>Range: 128-4096</p>
gab_isolate_time	<p>Maximum time to wait for isolated client</p> <p>Can be overridden at runtime</p> <p>See “About GAB run-time or dynamic tunable parameters” on page 240.</p>	<p>Default: 120000 msec (2 minutes)</p> <p>Range: 160000-240000 (in msec)</p>
gab_kill_ntries	<p>Number of times to attempt to kill client</p> <p>Can be overridden at runtime</p> <p>See “About GAB run-time or dynamic tunable parameters” on page 240.</p>	<p>Default: 5</p> <p>Range: 3-10</p>
gab_conn_wait	Maximum number of wait periods (as defined in the stable timeout parameter) before GAB disconnects the node from the cluster during cluster reconfiguration	<p>Default: 12</p> <p>Range: 1-256</p>

Table 5-1 GAB static tunable parameters (*continued*)

GAB parameter	Description	Values (default and range)
gab_ibuf_count	<p>Determines whether the GAB logging daemon is enabled or disabled</p> <p>The GAB logging daemon is enabled by default. To disable, change the value of <code>gab_ibuf_count</code> to 0.</p> <p>The disable login to the gab daemon while cluster is up and running with the <code>gabconfig -K</code> option. Use the <code>gabconfig</code> command to control value of this tunable.</p>	<p>Default: 8</p> <p>Range: 0-32</p>
gab_kstat_size	Number of system statistics to maintain in GAB	<p>Default: 60</p> <p>Range: 0 - 240</p>

About GAB run-time or dynamic tunable parameters

You can change the GAB dynamic tunable parameters while GAB is configured and while the cluster is running. The changes take effect immediately on running the `gabconfig` command. Note that some of these parameters also control how GAB behaves when it encounters a fault or a failure condition. Some of these conditions can trigger a PANIC which is aimed at preventing data corruption.

You can display the default values using the `gabconfig -l` command. To make changes to these values persistent across reboots, you can append the appropriate command options to the `/etc/gabtab` file along with any existing options. For example, you can add the `-k` option to an existing `/etc/gabtab` file that might read as follows:

```
gabconfig -c -n4
```

After adding the option, the `/etc/gabtab` file looks similar to the following:

```
gabconfig -c -n4 -k
```

Table 5-2 describes the GAB dynamic tunable parameters as seen with the `gabconfig -l` command, and specifies the command to modify them.

Table 5-2 GAB dynamic tunable parameters

GAB parameter	Description and command
Control port seed	<p>This option defines the minimum number of nodes that can form the cluster. This option controls the forming of the cluster. If the number of nodes in the cluster is less than the number specified in the <code>gabtab</code> file, then the cluster will not form. For example: if you type <code>gabconfig -c -n4</code>, then the cluster will not form until all four nodes join the cluster. If this option is enabled using the <code>gabconfig -x</code> command then the node will join the cluster even if the other nodes in the cluster are not yet part of the membership.</p> <p>Use the following command to set the number of nodes that can form the cluster:</p> <pre>gabconfig -n count</pre> <p>Use the following command to enable control port seed. Node can form the cluster without waiting for other nodes for membership:</p> <pre>gabconfig -x</pre>
Halt on process death	<p>Default: Disabled</p> <p>This option controls GAB's ability to halt (panic) the system on user process death. If <code>_had</code> and <code>_hashadow</code> are killed using <code>kill -9</code>, the system can potentially lose high availability. If you enable this option, then the GAB will PANIC the system on detecting the death of the client process. The default behavior is to disable this option.</p> <p>Use the following command to enable halt system on process death:</p> <pre>gabconfig -p</pre> <p>Use the following command to disable halt system on process death:</p> <pre>gabconfig -P</pre>

Table 5-2 GAB dynamic tunable parameters (*continued*)

GAB parameter	Description and command
Missed heartbeat halt	<p>Default: Disabled</p> <p>If this option is enabled then the system will panic on missing the first heartbeat from the VCS engine or the <code>vxconfigd</code> daemon in a CVM environment. The default option is to disable the immediate panic.</p> <p>This GAB option controls whether GAB can panic the node or not when the VCS engine or the <code>vxconfigd</code> daemon miss to heartbeat with GAB. If the VCS engine experiences a hang and is unable to heartbeat with GAB, then GAB will NOT PANIC the system immediately. GAB will first try to abort the process by sending SIGABRT (<code>kill_ntries</code> - default value 5 times) times after an interval of "iofence_timeout" (default value 15 seconds). If this fails, then GAB will wait for the "isolate_timeout" period which is controlled by a global tunable called <code>isolate_time</code> (default value 2 minutes). If the process is still alive, then GAB will PANIC the system.</p> <p>If this option is enabled GAB will immediately HALT the system in case of missed heartbeat from client.</p> <p>Use the following command to enable system halt when process heartbeat fails:</p> <pre>gabconfig -b</pre> <p>Use the following command to disable system halt when process heartbeat fails:</p> <pre>gabconfig -B</pre>

Table 5-2 GAB dynamic tunable parameters (*continued*)

GAB parameter	Description and command
Halt on rejoin	<p>Default: Disabled</p> <p>This option allows the user to configure the behavior of the VCS engine or any other user process when one or more nodes rejoin a cluster after a network partition. By default GAB will not PANIC the node running the VCS engine. GAB kills the userland process (the VCS engine or the vxconfigd process). This recycles the user port (port h in case of the VCS engine) and clears up messages with the old generation number programmatically. Restart of the process, if required, must be handled outside of GAB control, e.g., for hashadow process restarts _had.</p> <p>When GAB has kernel clients (such as fencing, VxVM, or VxFS), then the node will always PANIC when it rejoins the cluster after a network partition. The PANIC is mandatory since this is the only way GAB can clear ports and remove old messages.</p> <p>Use the following command to enable system halt on rejoin:</p> <pre data-bbox="581 829 736 855">gabconfig -j</pre> <p>Use the following command to disable system halt on rejoin:</p> <pre data-bbox="581 933 736 960">gabconfig -J</pre>
Keep on killing	<p>Default: Disabled</p> <p>If this option is enabled, then GAB prevents the system from PANICKING when the VCS engine or the vxconfigd process fail to heartbeat with GAB and GAB fails to kill the VCS engine or the vxconfigd process. GAB will try to continuously kill the VCS engine and will not panic if the kill fails.</p> <p>Repeat attempts to kill process if it does not die</p> <pre data-bbox="581 1260 736 1286">gabconfig -k</pre>

Table 5-2 GAB dynamic tunable parameters (*continued*)

GAB parameter	Description and command
Quorum flag	<p>Default: Disabled</p> <p>This is an option in GAB which allows a node to IOFENCE (resulting in a PANIC) if the new membership set is < 50% of the old membership set. This option is typically disabled and is used when integrating with other products</p> <p>Enable iofence quorum</p> <pre>gabconfig -q</pre> <p>Disable iofence quorum</p> <pre>gabconfig -d</pre>
GAB queue limit	<p>Default: Send queue limit: 128</p> <p>Default: Recv queue limit: 128</p> <p>GAB queue limit option controls the number of pending message before which GAB sets flow. Send queue limit controls the number of pending message in GAB send queue. Once GAB reaches this limit it will set flow control for the sender process of the GAB client. GAB receive queue limit controls the number of pending message in GAB receive queue before GAB send flow control for the receive side.</p> <p>Set the send queue limit to specified value</p> <pre>gabconfig -Q sendq:value</pre> <p>Set the receive queue limit to specified value</p> <pre>gabconfig -Q recvq:value</pre>
IOFENCE timeout	<p>Default: 15000(ms)</p> <p>This parameter specifies the timeout (in milliseconds) for which GAB will wait for the clients to respond to an IOFENCE message before taking next action. Based on the value of kill_ntries, GAB will attempt to kill client process by sending SIGABRT signal. If the client process is still registered after GAB attempted to kill client process for the value of kill_ntries times, GAB will halt the system after waiting for additional isolate_timeout value.</p> <p>Set the iofence timeout value to specified value in milliseconds.</p> <pre>gabconfig -f value</pre>

Table 5-2 GAB dynamic tunable parameters (*continued*)

GAB parameter	Description and command
Stable timeout	<p>Default: 5000(ms)</p> <p>Specifies the time GAB waits to reconfigure membership after the last report from LLT of a change in the state of local node connections for a given port. Any change in the state of connections will restart GAB waiting period.</p> <p>Set the stable timeout to specified value</p> <pre>gabconfig -t stable</pre>
Isolate timeout	<p>Default: 120000(ms)</p> <p>This tunable specifies the timeout value for which GAB will wait for client process to unregister in response to GAB sending SIGKILL signal. If the process still exists after isolate timeout GAB will halt the system</p> <pre>gabconfig -S isolate_time:value</pre>
Kill_ntries	<p>Default: 5</p> <p>This tunable specifies the number of attempts GAB will make to kill the process by sending SIGABRT signal.</p> <pre>gabconfig -S kill_ntries:value</pre>
Driver state	<p>This parameter shows whether GAB is configured. GAB may not have seeded and formed any membership yet.</p>
Partition arbitration	<p>This parameter shows whether GAB is asked to specifically ignore jeopardy.</p> <p>See the <code>gabconfig(1M)</code> manual page for details on the <code>-s</code> flag.</p>

About LLT tunable parameters

LLT provides various configuration and tunable parameters to modify and control the behavior of the LLT module. This section describes some of the LLT tunable parameters that can be changed at run-time and at LLT start-time.

See “[About Low Latency Transport \(LLT\)](#)” on page 25.

The tunable parameters are classified into two categories:

- LLT timer tunable parameters

See “[About LLT timer tunable parameters](#)” on page 246.

- LLT flow control tunable parameters

See “[About LLT flow control tunable parameters](#)” on page 250.

See “[Setting LLT timer tunable parameters](#)” on page 252.

About LLT timer tunable parameters

[Table 5-3](#) lists the LLT timer tunable parameters. The timer values are set in .01 sec units. The command `lltconfig -T query` can be used to display current timer values.

Table 5-3 LLT timer tunable parameters

LLT parameter	Description	Default	When to change	Dependency with other LLT tunable parameters
peerinact	LLT marks a link of a peer node as “inactive,” if it does not receive any packet on that link for this timer interval. Once a link is marked as “inactive,” LLT will not send any data on that link.	1600	<ul style="list-style-type: none"> ■ Change this value for delaying or speeding up node/link inactive notification mechanism as per client’s notification processing logic. ■ Increase the value for planned replacement of faulty network cable /switch. ■ In some circumstances, when the private networks links are very slow or the network traffic becomes very bursty, increase this value so as to avoid false notifications of peer death. Set the value to a high value for planned replacement of faulty network cable or faulty switch. 	The timer value should always be higher than the peertrouble timer value.

Table 5-3 LLT timer tunable parameters (*continued*)

LLT parameter	Description	Default	When to change	Dependency with other LLT tunable parameters
peertrouble	LLT marks a high-pri link of a peer node as "troubled", if it does not receive any packet on that link for this timer interval. Once a link is marked as "troubled", LLT will not send any data on that link till the link is up.	200	<ul style="list-style-type: none"> ■ In some circumstances, when the private networks links are very slow or nodes in the cluster are very busy, increase the value. ■ Increase the value for planned replacement of faulty network cable /faulty switch. 	This timer value should always be lower than peerinact timer value. Also, It should be close to its default value.
peertroublelo	LLT marks a low-pri link of a peer node as "troubled", if it does not receive any packet on that link for this timer interval. Once a link is marked as "troubled", LLT will not send any data on that link till the link is available.	400	<ul style="list-style-type: none"> ■ In some circumstances, when the private networks links are very slow or nodes in the cluster are very busy, increase the value. ■ Increase the value for planned replacement of faulty network cable /faulty switch. 	This timer value should always be lower than peerinact timer value. Also, It should be close to its default value.
heartbeat	LLT sends heartbeat packets repeatedly to peer nodes after every heartbeat timer interval on each highpri link.	50	In some circumstances, when the private networks links are very slow (or congested) or nodes in the cluster are very busy, increase the value.	This timer value should be lower than peertrouble timer value. Also, it should not be close to peertrouble timer value.
heartbeatlo	LLT sends heartbeat packets repeatedly to peer nodes after every heartbeatlo timer interval on each low pri link.	100	In some circumstances, when the networks links are very slow or nodes in the cluster are very busy, increase the value.	This timer value should be lower than peertroublelo timer value. Also, it should not be close to peertroublelo timer value.

Table 5-3 LLT timer tunable parameters (*continued*)

LLT parameter	Description	Default	When to change	Dependency with other LLT tunable parameters
timetoreqhb	If LLT does not receive any packet from the peer node on a particular link for "timetoreqhb" time period, it attempts to request heartbeats (sends 5 special heartbeat requests (hbreqs) to the peer node on the same link) from the peer node. If the peer node does not respond to the special heartbeat requests, LLT marks the link as "expired" for that peer node. The value can be set from the range of 0 to (peerinact -200). The value 0 disables the request heartbeat mechanism.	1400	Decrease the value of this tunable for speeding up node/link inactive notification mechanism as per client's notification processing logic. Disable the request heartbeat mechanism by setting the value of this timer to 0 for planned replacement of faulty network cable /switch. In some circumstances, when the private networks links are very slow or the network traffic becomes very bursty, don't change the value of this timer tunable.	This timer is set to 'peerinact - 200' automatically every time when the peerinact timer is changed.
reqhbtime	This value specifies the time interval between two successive special heartbeat requests. See the timetoreqhb parameter for more information on special heartbeat requests.	40	Symantec does not recommend to change this value	Not applicable
timetosendhb	LLT sends out of timer context heartbeats to keep the node alive when LLT timer does not run at regular interval. This option specifies the amount of time to wait before sending a heartbeat in case of timer not running. If this timer tunable is set to 0, the out of timer context heartbeating mechanism is disabled.	200	Disable the out of timer context heart-beating mechanism by setting the value of this timer to 0 for planned replacement of faulty network cable /switch. In some circumstances, when the private networks links are very slow or nodes in the cluster are very busy, increase the value	This timer value should not be more than peerinact timer value. Also, it should not be close to the peerinact timer value.

Table 5-3 LLT timer tunable parameters (*continued*)

LLT parameter	Description	Default	When to change	Dependency with other LLT tunable parameters
sendhbcap	This value specifies the maximum time for which LLT will send contiguous out of timer context heartbeats.	18000	Symantec does not recommend this value.	NA
oos	If the out-of-sequence timer has expired for a node, LLT sends an appropriate NAK to that node. LLT does not send a NAK as soon as it receives an oos packet. It waits for the oos timer value before sending the NAK.	10	Do not change this value for performance reasons. Lowering the value can result in unnecessary retransmissions/negative acknowledgement traffic. You can increase the value of oos if the round trip time is large in the cluster (for example, campus cluster).	Not applicable
retrans	LLT retransmits a packet if it does not receive its acknowledgement for this timer interval value.	10	Do not change this value. Lowering the value can result in unnecessary retransmissions. You can increase the value of retrans if the round trip time is large in the cluster (for example, campus cluster).	Not applicable
service	LLT calls its service routine (which delivers messages to LLT clients) after every service timer interval.	100	Do not change this value for performance reasons.	Not applicable
arp	LLT flushes stored address of peer nodes when this timer expires and relearns the addresses.	0	This feature is disabled by default.	Not applicable
arpreq	LLT sends an arp request when this timer expires to detect other peer nodes in the cluster.	3000	Do not change this value for performance reasons.	Not applicable

Table 5-3 LLT timer tunable parameters (*continued*)

LLT parameter	Description	Default	When to change	Dependency with other LLT tunable parameters
linkstable	This value specifies the amount of time to wait before LLT processes the link-down event for any link of the local node. LLT receives link-down events from the operating system when you enable the faster detection of link failure.	200	Increase this value in case of flaky links.	This timer value should not be more than peerinact timer value. Also, it should not be close to the peerinact timer value.

About LLT flow control tunable parameters

[Table 5-4](#) lists the LLT flow control tunable parameters. The flow control values are set in number of packets. The command `lltconfig -F query` can be used to display current flow control settings.

Table 5-4 LLT flow control tunable parameters

LLT parameter	Description	Default	When to change	Dependency with other LLT tunable parameters
highwater	When the number of packets in transmit queue for a node reaches highwater, LLT is flow controlled.	200	If a client generates data in bursty manner, increase this value to match the incoming data rate. Note that increasing the value means more memory consumption so set an appropriate value to avoid wasting memory unnecessarily. Lowering the value can result in unnecessary flow controlling the client.	This flow control value should always be higher than the lowwater flow control value.
lowwater	When LLT has flow controlled the client, it will not start accepting packets again till the number of packets in the port transmit queue for a node drops to lowwater.	100	Symantec does not recommend to change this tunable.	This flow control value should be lower than the highwater flow control value. The value should not be close the highwater flow control value.

Table 5-4 LLT flow control tunable parameters (*continued*)

LLT parameter	Description	Default	When to change	Dependency with other LLT tunable parameters
rpothighwater	When the number of packets in the receive queue for a port reaches highwater, LLT is flow controlled.	200	If a client generates data in bursty manner, increase this value to match the incoming data rate. Note that increasing the value means more memory consumption so set an appropriate value to avoid wasting memory unnecessarily. Lowering the value can result in unnecessary flow controlling the client on peer node.	This flow control value should always be higher than the rportlowwater flow control value.
rportlowwater	When LLT has flow controlled the client on peer node, it will not start accepting packets for that client again till the number of packets in the port receive queue for the port drops to rportlowwater.	100	Symantec does not recommend to change this tunable.	This flow control value should be lower than the rpothighwater flow control value. The value should not be close the rpothighwater flow control value.
window	This is the maximum number of un-ACKed packets LLT will put in flight.	50	Change the value as per the private networks speed. Lowering the value irrespective of network speed may result in unnecessary retransmission of out of window sequence packets.	This flow control value should not be higher than the difference between the highwater flow control value and the lowwater flow control value. The value of this parameter (window) should be aligned with the value of the bandwidth delay product.
linkburst	It represents the number of back-to-back packets that LLT sends on a link before the next link is chosen.	32	For performance reasons, its value should be either 0 or at least 32.	This flow control value should not be higher than the difference between the highwater flow control value and the lowwater flow control value.

Table 5-4 LLT flow control tunable parameters (*continued*)

LLT parameter	Description	Default	When to change	Dependency with other LLT tunable parameters
ackval	LLT sends acknowledgement of a packet by piggybacking an ACK packet on the next outbound data packet to the sender node. If there are no data packets on which to piggyback the ACK packet, LLT waits for ackval number of packets before sending an explicit ACK to the sender.	10	Do not change this value for performance reasons. Increasing the value can result in unnecessary retransmissions.	Not applicable
sws	To avoid Silly Window Syndrome, LLT transmits more packets only when the count of un-acked packet goes to below of this tunable value.	40	For performance reason, its value should be changed whenever the value of the window tunable is changed as per the formula given below: sws = window *4/5.	Its value should be lower than that of window. Its value should be close to the value of window tunable.
largepktlen	When LLT has packets to delivers to multiple ports, LLT delivers one large packet or up to five small packets to a port at a time. This parameter specifies the size of the large packet.	1024	Symantec does not recommend to change this tunable.	Not applicable

Setting LLT timer tunable parameters

You can set the LLT tunable parameters either with the `lltconfig` command or in the `/etc/llttab` file. You can use the `lltconfig` command to change a parameter on the local node at run time. Symantec recommends you run the command on all the nodes in the cluster to change the values of the parameters. To set an LLT parameter across system reboots, you must include the parameter definition in the `/etc/llttab` file. Default values of the parameters are taken if nothing is specified in `/etc/llttab`. The parameters values specified in the `/etc/llttab` file come into effect at LLT start-time only. Symantec recommends that you specify the same definition of the tunable parameters in the `/etc/llttab` file of each node.

To get and set a timer tunable:

- To get the current list of timer tunable parameters using `lltconfig` command:

```
# lltdconfig -T query
```

- To set a timer tunable parameter using the `lltdconfig` command:

```
# lltdconfig -T timer tunable:value
```

- To set a timer tunable parameter in the `/etc/lltdtab` file:

```
set-timer timer tunable:value
```

To get and set a flow control tunable

- To get the current list of flow control tunable parameters using `lltdconfig` command:

```
# lltdconfig -F query
```

- To set a flow control tunable parameter using the `lltdconfig` command:

```
# lltdconfig -F flowcontrol tunable:value
```

- To set a flow control tunable parameter in the `/etc/lltdtab` file:

```
set-flow flowcontrol tunable:value
```

See the `lltdconfig(1M)` and `lltdtab(1M)` manual pages.

About LMX tunable parameters

The section describes the LMX tunable parameters and how to reconfigure the LMX module.

LMX tunable parameters

[Table 5-5](#) describes the LMX driver tunable parameters.

Table 5-5 LMX Tunable parameters

LMX parameter	Default value	Maximum value	Description
lmx_minor_max	8192	65535	Specifies the maximum number of contexts system-wide. Each Oracle process typically has two LMX contexts. "Contexts" and "minors" are used interchangeably in the documentation; "context" is an Oracle-specific term to specify the value in the lmx.conf file.
lmx_port_max	8192	65535	Specifies the number of communication endpoints for transferring messages from the sender to receiver in a uni-directional manner.
lmx_buffer_max	8192	65535	Specifies the number of addressable regions in memory to copy LMX data.

If you see the message "no minors available" on one node, add a configuration parameter that increases the value for the maximum number of contexts.

Note: The error message may contain the term "minors," but you must use the term "contexts" when changing the parameter value.

Warning: Increasing the number of contexts on a specific system has some impact on the resources of that system.

Reconfiguring the LMX module

This section discusses how to reconfigure the LMX module on the node. For the parameter changes to take effect, you must reconfigure the LMX module.

To reconfigure the LMX module

- 1 Configure the tunable parameter.

```
# /usr/sbin/kctune tunable=value
```

For example:

```
# /usr/sbin/kctune lmx_minor_max=16384
```

- 2 If you use Oracle RAC 10g or Oracle RAC 11g , stop Oracle Clusterware (if Oracle Clusterware is not under VCS control) and verify that Oracle Clusterware is stopped. You must also stop applications which are not under VCS control.
- 3 Unmount the CFS mounts (if mounts are not under VCS control).
- 4 Stop VCS by entering the following command:

```
# /opt/VRTSvcs/bin/hastop -local
```

- 5 Check that this node is registered at gab ports a, b, d, and o only. Ports f, h, u, v, y, and w should not be seen on this node.

```
# gabconfig -a
```

```
GAB Port Memberships
```

```
=====
```

```
Port a gen ada401 membership 0123
```

```
Port b gen ada40d membership 0123
```

```
Port d gen ada409 membership 0123
```

```
Port o gen ada406 membership 0123
```

- 6 Restart the node by entering the following command:

```
sys1> # /usr/sbin/shutdown -r
```

CFS mounts controlled by VCS are automatically remounted, but you must manually remount CFS mounts which are not under VCS control.

Applications which are outside of VCS control must be manually restarted.

About VXFEN tunable parameters

The section describes the VXFEN tunable parameters and how to reconfigure the VXFEN module.

[Table 5-6](#) describes the tunable parameters for the VXFEN driver.

Table 5-6 VXFEN tunable parameters

vxfen Parameter	Description and Values: Default, Minimum, and Maximum
vxfen_debug_sz	<p>Size of debug log in bytes</p> <ul style="list-style-type: none"> ■ Values Default: 131072 (128 KB) Minimum: 65536 (64 KB) Maximum: 524288 (512 KB)
vxfen_max_delay	<p>Specifies the maximum number of seconds that the smaller sub-cluster waits before racing with larger sub-clusters for control of the coordinator disks when a network partition occurs.</p> <p>This value must be greater than the vxfen_min_delay value.</p> <ul style="list-style-type: none"> ■ Values Default: 60 Minimum: 1 Maximum: 600
vxfen_min_delay	<p>Specifies the minimum number of seconds that the smaller sub-cluster waits before racing with larger sub-clusters for control of the coordinator disks when a network partition occurs.</p> <p>This value must be smaller than or equal to the vxfen_max_delay value.</p> <ul style="list-style-type: none"> ■ Values Default: 1 Minimum: 1 Maximum: 600
vxfen_vxfnd_tmt	<p>Specifies the time in seconds that the I/O fencing driver Vxfen waits for the I/O fencing daemon VXFEND to return after completing a given task.</p> <ul style="list-style-type: none"> ■ Values Default: 60 Minimum: 10 Maximum: 600

Table 5-6 VXFEN tunable parameters (*continued*)

vxfen Parameter	Description and Values: Default, Minimum, and Maximum
panic_timeout_offst	<p>Specifies the time in seconds based on which the I/O fencing driver VxFEN computes the delay to pass to the GAB module to wait until fencing completes its arbitration before GAB implements its decision in the event of a split-brain. You can set this parameter in the <code>vxfenmode</code> file and use the <code>vxfenadm</code> command to check the value. Depending on the <code>vxfen_mode</code>, the GAB delay is calculated as follows:</p> <ul style="list-style-type: none"> ■ For scsi3 mode: $1000 * (\text{panic_timeout_offst} + \text{vxfen_max_delay})$ ■ For customized mode: $1000 * (\text{panic_timeout_offst} + \max(\text{vxfen_vxwnd_tmt}, \text{vxfen_loser_exit_delay}))$ ■ Default: 10

In the event of a network partition, the smaller sub-cluster delays before racing for the coordinator disks. The time delay allows a larger sub-cluster to win the race for the coordinator disks. The `vxfen_max_delay` and `vxfen_min_delay` parameters define the delay in seconds.

Configuring the VXFEN module parameters

After adjusting the tunable kernel driver parameters, you must reconfigure the VXFEN module for the parameter changes to take effect.

The following example procedure changes the value of the `vxfen_min_delay` parameter.

To configure the VxFEN parameters and reconfigure the VxFEN module

- 1 Configure the tunable parameter.

```
# /usr/sbin/kctune tunable=value
```

For example:

```
# /usr/sbin/kctune vxfen_min_delay=100
```

- 2 If you use Oracle 10g or Oracle 11g, stop CRS (if CRS is not under VCS control) and verify that CRS is stopped.

- 3 Unmount CFS mounts (if mounts are not under VCS control).

Determine the file systems to unmount by checking the `/etc/mnttab` file.

```
# mount | grep vxfs | grep cluster
```

To unmount the mount points listed in the output, enter:

```
# umount mount_point
```

- 4 Stop VCS.

```
# hastop -local
```

- 5 Check that this node is registered at gab ports a, b, d, and o only.

Ports f, h, u, v, w, and y should not be seen on this node.

```
# gabconfig -a
```

```
GAB Port Memberships
```

```
=====
Port a gen ada401 membership 0123
Port b gen ada40d membership 0123
Port d gen ada409 membership 0123
Port o gen ada406 membership 0123
```

- 6 Reboot the node.

```
# /usr/sbin/shutdown -r
```

Tuning guidelines for campus clusters

An important consideration while tuning an SF Oracle RAC campus cluster is setting the LLT peerinact time. Follow the guidelines below to determine the optimum value of peerinact time:

- Calculate the roundtrip time using `lltping (1M)`.
- Evaluate LLT heartbeat time as half of the round trip time.
- Set the LLT peer trouble time as 2-4 times the heartbeat time.
- LLT peerinact time should be set to be more than 4 times the heart beat time.

Reference

- [Appendix A. List of SF Oracle RAC health checks](#)
- [Appendix B. Error messages](#)

List of SF Oracle RAC health checks

This appendix includes the following topics:

- [LLT health checks](#)
- [LMX health checks in SF Oracle RAC clusters](#)
- [I/O fencing health checks](#)
- [PrivNIC health checks in SF Oracle RAC clusters](#)
- [Oracle Clusterware health checks in SF Oracle RAC clusters](#)
- [CVM, CFS, and ODM health checks in SF Oracle RAC clusters](#)

LLT health checks

This section lists the health checks performed for LLT, the messages displayed for each check, and a brief description of the check.

Note: Warning messages indicate issues in the components or the general health of the cluster. Follow the troubleshooting recommendations to resolve the issues.

See [“Troubleshooting LLT health check warning messages”](#) on page 195.

[Table A-1](#) lists the health checks performed for LLT.

Table A-1 List of health checks for LLT

List of health checks	Message	Description
LLT timer subsystem scheduling check	Warning: OS timer is not called for <i>num</i> seconds	<p>Checks whether the LLT module runs in accordance with the scheduled operating system timer.</p> <p>The message indicates that the operating system timer is not called for the specified interval.</p> <p>The parameter <code>timer_threshold</code> contains the optimum threshold for this check.</p>
OS memory availability check for the packet transmission	Warning: Kernel failed to allocate memory <i>num</i> time(s)	<p>Checks whether the kernel has allocated sufficient memory to LLT for cluster communication.</p> <p>The message indicates that the kernel attempts at allocating the requisite memory has failed (<i>num</i>) times.</p> <p>The parameter <code>no_mem_to_xmit_allow</code> contains the optimum threshold for this check.</p>
Flow control status monitoring	Flow-control occurred <i>num</i> time(s) and back-enabled <i>num</i> time(s) on port <i>port number</i> for node <i>node number</i>	<p>Checks whether LLT has sufficient bandwidth to accept incoming data packets and transmit the data packets. Flow control also depends on the peer nodes' ability to accept incoming data traffic.</p> <p>The message indicates that packet transmission and reception was controlled (<i>num</i>) and normalized (<i>num</i>) times.</p> <p>The parameter <code>max_canput</code> contains the optimum threshold for this check.</p>

Table A-1 List of health checks for LLT (*continued*)

List of health checks	Message	Description
Flaky link monitoring	Warning: Connectivity with node <i>node id</i> on link <i>link id</i> is flaky <i>num</i> time(s). The distribution is: (0-4 s) <num> (4-8 s) <num> (8-12 s) <num> (12-16 s) <num> (>=16 s)	Checks whether the private interconnects are stable. The message indicates that connectivity (<i>link</i>) with the peer node (<i>node id</i>) is monitored (<i>num</i>) times within a stipulated duration (for example, 0-4 seconds). The parameter <code>flakylink_allow</code> contains the optimum threshold for this check.
Node status check	One or more link connectivity with peer node(s) <i>node name</i> is in trouble.	Checks the current node membership. The message indicates connectivity issues with the peer node (<i>node name</i>). This check does not require a threshold.
Link status check	Link connectivity with <i>node name</i> is on only one link. Symantec recommends configuring a minimum of 2 links.	Checks the number of private interconnects that are currently active. The message indicates that only one private interconnect exists or is operational between the local node and the peer node <i>node id</i> . This check does not require a threshold.
Connectivity check	Only one link is configured under LLT. Symantec recommends configuring a minimum of 2 links.	Checks the number of private interconnects that were configured under LLT during configuration. The message indicates that only one private interconnect exists or is configured under LLT. This check does not require a threshold.

Table A-1 List of health checks for LLT (*continued*)

List of health checks	Message	Description
LLT packet related checks	<p>Retransmitted % percentage of total transmitted packets.</p> <p>Sent % percentage of total transmitted packet when no link is up.</p> <p>% percentage of total received packets are with bad checksum.</p> <p>% percentage of total received packets are out of window.</p> <p>% percentage of total received packets are misaligned.</p>	<p>Checks whether the data packets transmitted by LLT reach peer nodes without any error. If there is an error in packet transmission, it indicates an error in the private interconnects.</p> <p>The message indicates the percentage of data packets transmitted or received and the associated transmission errors, such as bad checksum and failed link.</p> <p>The following parameters contain the optimum thresholds for this check: <code>retrans_pct</code>, <code>send_nolinkup_pct</code>, <code>recv_oow_pct</code>, <code>recv_misaligned_pct</code>, <code>recv_badcksum_pct</code></p>
DLPI layer related checks	<p>% percentage of total received packets are with DLPI error.</p>	<p>Checks whether the data link layer (DLP) is causing errors in packet transmission.</p> <p>The message indicates that some of the received packets (%) contain DLPI error.</p> <p>You can set a desired percentage value in the configuration file.</p> <p>The parameter <code>recv_dlpierror_pct</code> contains the optimum threshold for this check.</p>

Table A-1 List of health checks for LLT (*continued*)

List of health checks	Message	Description
Traffic distribution over the links	Traffic distribution over links: %% Send data on linknum percentage %% Recv data on linknum percentage	Checks the distribution of traffic over all the links configured under LLT. The message displays the percentage of data (%) sent and recd on a particular link (<i>num</i>) This check does not require a threshold.
LLT Ports status check	% per of total transmitted packets are with large xmit latency (>16ms) for port port id %per received packets are with large recv latency (>16ms) for port port id.	Checks the latency period for transmitting or receiving packets. The message indicates that some percentage (%) of the transmitted/received packets exceed the stipulated latency time. The following parameters contain the optimum thresholds for this check: hirecvlatencycnt_pct, hixmitlatencycnt_pct.
System load monitoring	Load Information: Average : num, num, num	Monitors the system workload at the stipulated periodicity (1 second, 5 seconds, 15 seconds) This check does not require a threshold.

LMX health checks in SF Oracle RAC clusters

This section lists the health checks performed for LMX, the messages displayed for each check, and a brief description of the check.

Note: Warning messages indicate issues in the components or the general health of the cluster.

For recommendations on resolving the issues, see the troubleshooting chapter in this document.

Table A-2 lists the health checks performed for LMX.

Table A-2 List of health checks for LMX

List of health checks	Message	Description
Checks related to the working status of LMX/VCSMM driver	<p>LMX is not running. This warning is not applicable for Oracle 11g running cluster.</p> <p>VCSMM is not running.</p>	<p>Checks whether VCSMM/LMX is running. This warning is valid only for clusters running Oracle RAC 10g.</p> <p>Note: This check is performed only if the database cache fusion traffic occurs over VCSIPC.</p>
Checks related to confirm linking of Oracle binary with the LMX/VCSMM library.	<p>Oracle is not linked to the Symantec LMX library. This warning is not applicable for Oracle 11g running cluster.</p> <p>Oracle is linked to Symantec LMX Library, but LMX is not running.</p> <p>Oracle is not linked to Symantec VCSMM library.</p> <p>Oracle is linked to Symantec VCSMM Library, but VCSMM is not running.</p>	<p>Checks whether the Oracle RAC binary is linked with the LMX/VCSMM library. This warning is valid only for clusters running Oracle RAC 10g.</p> <p>Note: This check is performed only if the database cache fusion traffic occurs over VCSIPC.</p>
Oracle traffic on LMX port related information	<p>Traffic over LMX <i>numnum num</i></p>	<p>Displays the Oracle RAC traffic over LMX if the database cache fusion uses VCSIPC for inter-process communication.</p> <p>The message displays the traffic over LMX at the time that the script is invoked.</p>

I/O fencing health checks

This section lists the health checks performed for I/O fencing, the messages displayed for each check, and a brief description of the check.

Note: Warning messages indicate issues in the components or the general health of the cluster.

For recommendations on resolving the issues, see the troubleshooting chapter in this document.

[Table A-3](#) lists the health checks performed for I/O fencing.

Table A-3 List of health checks for I/O fencing

List of health checks	Message	Description
Checks related to the working status of VxFEN	VxFEN is not configured in SCSI3 mode. VxFEN is not running on all nodes in the cluster.	Checks whether Symantec fencing is configured in SCSI3 mode and running on all nodes in the cluster.
Checks related to coordinator disk health	VxFEN is running with only one coordinator disk. Loss of this disk will prevent cluster reconfiguration on loss of a node. Symantec recommends configuring a minimum of 3 coordinator disks. VxFEN is running with even number of coordinator disks. There must be odd number of coordinator disks. Replace the disk <i>disk name</i> using OCDR procedure.	Checks how many coordinator disks are configured for Symantec fencing and the status of the disks.
Verify the keys on the coordinator disks	The coordinator disk (<i>disk name</i>) does not have the required key for the local node.	Checks whether the coordinator disks configured under Symantec fencing has the key of the current node.

Table A-3 List of health checks for I/O fencing (*continued*)

List of health checks	Message	Description
Verify SCSI3 write-exclusive reservation on shared disk	SCSI3 write-exclusive reservation is missing on shared disk <i>(disk_name)</i>	Checks if the shared disk is accessible to the node for write operations.

PrivNIC health checks in SF Oracle RAC clusters

This section lists the health checks performed for PrivNIC, the messages displayed for each check, and a brief description of the check.

Note: Warning messages indicate issues in the components or the general health of the cluster.

For recommendations on resolving the issues, see the troubleshooting chapter in this document.

[Table A-4](#) lists the health checks performed for PrivNIC.

Table A-4 List of health checks for PrivNIC

List of health checks	Message	Description
Checks related to the working status of VCS	VCS is not running on the local system.	Checks whether VCS is running or not.
Checks related to the status of the PrivNIC resources	If the PrivNIC resource is not online: The PrivNIC resource <i>resource name</i> is not online.	Checks whether the PrivNIC resources configured under VCS are online.
Compare the NICs used by PrivNIC with the NICs configured under LLT	For PrivNIC resources: Mismatch between LLT links <i>llt nics</i> and PrivNIC links <i>private nics</i> .	Checks whether the NICs used by the PrivNIC resource have the same interface (<i>private nics</i>) as those configured as LLT links (<i>llt nics</i>).

Table A-4 List of health checks for PrivNIC (*continued*)

List of health checks	Message	Description
Cross check NICs used by PrivNIC with the NICs used by Oracle Clusterware.	For PrivNIC resources: Mismatch between Oracle Clusterware links <i>crs_nics</i> and PrivNIC links <i>private_nics</i> .	Checks whether the private interconnect used for Oracle Clusterware is the same as the NIC configured under the PrivNIC resource.

Oracle Clusterware health checks in SF Oracle RAC clusters

This section lists the health checks performed for Oracle Clusterware, the messages displayed for each check, and a brief description of the check.

Note: Warning messages indicate issues in the components or the general health of the cluster.

For recommendations on resolving the issues, see the troubleshooting chapter in this document.

[Table A-5](#) lists the health checks performed for the Oracle Clusterware module.

Table A-5 List of health checks for the Oracle Clusterware module

List of health checks	Message	Description
Verify the working status of CRS	Oracle Clusterware is not running.	Checks whether Oracle Clusterware is up and running.
Verify the working status of CSSD resource	No CSSD resource is configured under VCS. The CSSD resource <i>name</i> is not running.	Checks whether the CSSD resource is configured under VCS and if the resource is online.
Compare the NICs used by CRS with the NICs configured under LLT	Mismatch between LLT links <i>llt_nic1</i> , <i>llt_nic2</i> and Oracle Clusterware links <i>crs_nic1</i> , <i>crs_nic2</i> .	Checks whether the private interconnects used by Oracle Clusterware are the same as the LLT links (<i>llt_nics</i>).

Table A-5 List of health checks for the Oracle Clusterware module (*continued*)

List of health checks	Message	Description
Compare the NICs used by CRS with the NICs monitored by PrivNIC	Mismatch between Oracle Clusterware links <i>crs nics</i> and PrivNIC links <i>private nics</i> .	Checks whether the private interconnects (<i>crs nics</i>) used by Oracle Clusterware are monitored by the PrivNIC resource (<i>private nics</i>).
Compare the nodes configured for CRS with the nodes listed by llc commands.	Mismatch between CRS nodes <i>crs nodes</i> and LLT nodes <i>llt nodes</i> .	Checks whether the host names configured during the Oracle Clusterware installation match the list of host names displayed on running the LLT command: <code>lltstat -nvv</code> .

CVM, CFS, and ODM health checks in SF Oracle RAC clusters

This section lists the health checks performed for CVM, CFS, and ODM, the messages displayed for each check, and a brief description of the check.

Note: Warning messages indicate issues in the components or the general health of the cluster.

For recommendations on resolving the issues, see the troubleshooting chapter in this document.

[Table A-6](#) lists the health checks performed for CVM, CFS, and ODM modules.

Table A-6 List of health checks for CVM, CFS, and ODM

List of health checks	Message	Description
Verify CVM status	CVM is not running	Checks whether CVM is running in the cluster.
Verify CFS status	CFS is not running	Checks whether CFS is running in the cluster.
Verify ODM status	ODM is not running	Checks whether ODM is running in the cluster.

Error messages

This appendix includes the following topics:

- [About error messages](#)
- [LMX error messages in SF Oracle RAC](#)
- [VxVM error messages](#)
- [VXFEN driver error messages](#)

About error messages

Error messages can be generated by the following software modules:

- LLT Multiplexer (LMX)
- Veritas Volume Manager (VxVM)
- Veritas Fencing (VXFEN) driver

LMX error messages in SF Oracle RAC

There are two types of LMX error messages: critical and non-critical.

Gather information about the systems and system configurations prior to contacting Symantec support personnel for assistance with error messages. This information will assist Symantec support personnel in identifying and resolving the error.

See [“Gathering information from an SF Oracle RAC cluster for support analysis”](#) on page 180.

LMX critical error messages in SF Oracle RAC

The messages in [Table B-1](#) report critical errors seen when the system runs out of memory, when LMX is unable to communicate with LLT, or when you are unable to load or unload LMX.

[Table B-1](#) lists the critical LMX kernel module error messages.

Table B-1 LMX critical error messages

Message ID	LMX Message
00001	lmxload packet header size incorrect (number)
00002	lmxload invalid lmx_llt_port number
00003	lmxload context memory alloc failed
00004	lmxload port memory alloc failed
00005	lmxload buffer memory alloc failed
00006	lmxload node memory alloc failed
00007	lmxload msgbuf memory alloc failed
00008	lmxload tmp msgbuf memory alloc failed
00009	lmxunload node number conngrp not NULL
00010	lmxopen return, minor non-zero
00011	lmxopen return, no minors available
00012	lmxconnect lmxlltopen(1) err= number
00013	lmxconnect new connection memory alloc failed
00014	lmxconnect kernel request memory alloc failed
00015	lmxconnect mblk memory alloc failed
00016	lmxconnect conn group memory alloc failed
00017	lmxlltunregister: LLT unregister port number failed err= number
00018	lmxload contexts number > number, max contexts = system limit = number
00019	lmxload ports number > number, max ports = system limit = number

Table B-1 LMX critical error messages (*continued*)

Message ID	LMX Message
00020	lmxload buffers number > number, max buffers = system limit = number
00021	lmxload msgbuf number > number, max msgbuf size = system limit = number

LMX non-critical error messages in SF Oracle RAC

If the message displays in [Table B-2](#) creates errors while running an Oracle application, use the `lmxconfig` command to turn off the display. For example:

```
# lmxconfig -e 0
```

To re-enable message displays, type:

```
# lmxconfig -e 1
```

[Table B-2](#) contains LMX error messages that may appear during run-time.

Table B-2 LMX non-critical error messages

Message ID	LMX Message
06001	lmxreqlink duplicate kreq= 0xaddress, req= 0xaddress
06002	lmxreqlink duplicate ureq= 0xaddress kr1= 0xaddress, kr2= 0xaddress req type = number
06003	lmxrequnlink not found kreq= 0xaddress from= number
06004	lmxrequnlink_l not found kreq= 0xaddress from= number
06101	lmxpollreq not in doneq CONN kreq= 0xaddress
06201	lmxnewcontext lltinit fail err= number
06202	lmxnewcontext llregister fail err= number
06301	lmxrecvport port not found unode= number node= number ctx= number
06302	lmxrecvport port not found (no port) ctx= number
06303	lmxrecvport port not found ugen= number gen= number ctx= number
06304	lmxrecvport dup request detected

Table B-2 LMX non-critical error messages (*continued*)

Message ID	LMX Message
06401	lmxinitport out of ports
06501	lmxsendport lltsend node= number err= number
06601	lmxinitbuf out of buffers
06602	lmxinitbuf fail ctx= number ret= number
06701	lmxsendbuf lltsend node= number err= number
06801	lmxconfig insufficient privilege, uid= number
06901	lmxlltnodestat: LLT getnodeinfo failed err= number

VxVM error messages

[Table B-3](#) contains VxVM error messages that are related to I/O fencing.

Table B-3 VxVM error messages for I/O fencing

Message	Explanation
<code>vold_pgr_register(disk_path): failed to open the vxfen device. Please make sure that the vxfen driver is installed and configured.</code>	The vxfen driver is not configured. Follow the instructions to set up these disks and start I/O fencing. You can then clear the faulted resources and bring the service groups online.
<code>vold_pgr_register(disk_path): Probably incompatible vxfen driver.</code>	Incompatible versions of VxVM and the vxfen driver are installed on the system. Install the proper version of SF Oracle RAC.

VXFEN driver error messages

[Table B-4](#) contains VXFEN driver error messages. In addition to VXFEN driver error messages, informational messages can also be displayed.

See [“VXFEN driver informational message”](#) on page 275.

See [“Node ejection informational messages”](#) on page 275.

Table B-4 VXFEN driver error messages

Message	Explanation
Unable to register with coordinator disk with serial number: xxxx	This message appears when the vxfen driver is unable to register with one of the coordinator disks. The serial number of the coordinator disk that failed is displayed.
Unable to register with a majority of the coordinator disks. Dropping out of cluster.	<p>This message appears when the vxfen driver is unable to register with a majority of the coordinator disks. The problems with the coordinator disks must be cleared before fencing can be enabled.</p> <p>This message is preceded with the message "VXFEN: Unable to register with coordinator disk with serial number xxxx."</p>

VXFEN driver informational message

The following informational message appears when a node is ejected from the cluster to prevent data corruption when a split-brain occurs.

```
VXFEN CRITICAL V-11-1-20 Local cluster node ejected from cluster  
to prevent potential data corruption
```

Node ejection informational messages

Informational messages may appear on the console of one of the cluster nodes when a node is ejected from a disk or LUN.

These informational messages can be ignored.

Glossary

Agent	A process that starts, stops, and monitors all configured resources of a type, and reports their status to VCS.
Authentication Broker	The Veritas Security Services component that serves, one level beneath the root broker, as an intermediate registration authority and a certification authority. The authentication broker can authenticate clients, such as users or services, and grant them a certificate that will become part of the Veritas credential. An authentication broker cannot, however, authenticate other brokers. That task must be performed by the root broker.
Cluster	A cluster is one or more computers that are linked together for the purpose of multiprocessing and high availability. The term is used synonymously with VCS cluster, meaning one or more computers that are part of the same GAB membership.
CVM (cluster volume manager)	The cluster functionality of Veritas Volume Manager.
Disaster Recovery	Administrators with clusters in physically disparate areas can set the policy for migrating applications from one location to another if clusters in one geographic area become unavailable due to an unforeseen event. Disaster recovery requires heartbeating and replication.
disk array	A collection of disks logically arranged into an object. Arrays tend to provide benefits such as redundancy or improved performance.
DMP (Dynamic Multi-Pathing)	A feature designed to provide greater reliability and performance by using path failover and load balancing for multiported disk arrays connected to host systems through multiple paths. DMP detects the various paths to a disk using a mechanism that is specific to each supported array type. DMP can also differentiate between different enclosures of a supported array type that are connected to the same host system.
SmartTier	A feature with which administrators of multi-volume VxFS file systems can manage the placement of files on individual volumes in a volume set by defining placement policies that control both initial file location and the circumstances under which existing files are relocated. These placement policies cause the files to which they apply to be created and extended on specific subsets of a file system's volume set, known as placement classes. The files are relocated to volumes in other placement

classes when they meet specified naming, timing, access rate, and storage capacity-related conditions.

Failover	A failover occurs when a service group faults and is migrated to another system.
GAB (Group Atomic Broadcast)	A communication mechanism of the VCS engine that manages cluster membership, monitors heartbeat communication, and distributes information throughout the cluster.
HA (high availability)	The concept of configuring the product to be highly available against system failure on a clustered network using Veritas Cluster Server (VCS).
HAD (High Availability Daemon)	The core VCS process that runs on each system. The HAD process maintains and communicates information about the resources running on the local system and receives information about resources running on other systems in the cluster.
IP address	<p>An identifier for a computer or other device on a TCP/IP network, written as four eight-bit numbers separated by periods. Messages and other data are routed on the network according to their destination IP addresses.</p> <p>See also virtual IP address</p>
Jeopardy	A node is in jeopardy when it is missing one of the two required heartbeat connections. When a node is running with one heartbeat only (in jeopardy), VCS does not restart the applications on a new node. This action of disabling failover is a safety mechanism that prevents data corruption.
latency	For file systems, this typically refers to the amount of time it takes a given file system operation to return to the user.
LLT (Low Latency Transport)	A communication mechanism of the VCS engine that provides kernel-to-kernel communications and monitors network communications.
logical volume	<p>A simple volume that resides on an extended partition on a basic disk and is limited to the space within the extended partitions. A logical volume can be formatted and assigned a drive letter, and it can be subdivided into logical drives.</p> <p>See also LUN</p>
LUN	A LUN, or logical unit, can either correspond to a single physical disk, or to a collection of disks that are exported as a single logical entity, or virtual disk, by a device driver or by an intelligent disk array's hardware. VxVM and other software modules may be capable of automatically discovering the special characteristics of LUNs, or you can use disk tags to define new storage attributes. Disk tags are administered by using the <code>vxdisk</code> command or the graphical user interface.
main.cf	The file in which the cluster configuration is stored.
mirroring	A form of storage redundancy in which two or more identical copies of data are maintained on separate volumes. (Each duplicate copy is known as a mirror.) Also RAID Level 1.

Node	The physical host or system on which applications and service groups reside. When systems are linked by VCS, they become nodes in a cluster.
resources	Individual components that work together to provide application services to the public network. A resource may be a physical component such as a disk group or network interface card, a software component such as a database server or a Web server, or a configuration component such as an IP address or mounted file system.
Resource Dependency	A dependency between resources is indicated by the keyword "requires" between two resource names. This indicates the second resource (the child) must be online before the first resource (the parent) can be brought online. Conversely, the parent must be offline before the child can be taken offline. Also, faults of the children are propagated to the parent.
Resource Types	Each resource in a cluster is identified by a unique name and classified according to its type. VCS includes a set of pre-defined resource types for storage, networking, and application services.
root broker	The first authentication broker, which has a self-signed certificate. The root broker has a single private domain that holds only the names of brokers that shall be considered valid.
SAN (storage area network)	A networking paradigm that provides easily reconfigurable connectivity between any subset of computers, disk storage and interconnecting hardware such as switches, hubs and bridges.
Service Group	A service group is a collection of resources working together to provide application services to clients. It typically includes multiple resources, hardware- and software-based, working together to provide a single service.
Service Group Dependency	A service group dependency provides a mechanism by which two service groups can be linked by a dependency rule, similar to the way resources are linked.
Shared Storage	Storage devices that are connected to and used by two or more systems.
shared volume	A volume that belongs to a shared disk group and is open on more than one node at the same time.
SNMP Notification	Simple Network Management Protocol (SNMP) developed to manage nodes on an IP network.
State	The current activity status of a resource, group or system. Resource states are given relative to both systems.
Storage Checkpoint	A facility that provides a consistent and stable view of a file system or database image and keeps track of modified data blocks since the last Storage Checkpoint.
System	The physical system on which applications and service groups reside. When a system is linked by VCS, it becomes a node in a cluster. See Node

types.cf	A file that describes standard resource types to the VCS engine; specifically, the data required to control a specific resource.
VCS (Veritas Cluster Server)	An open systems clustering solution designed to eliminate planned and unplanned downtime, simplify server consolidation, and allow the effective management of a wide range of applications in multiplatform environments.
Virtual IP Address	A unique IP address associated with the cluster. It may be brought up on any system in the cluster, along with the other resources of the service group. This address, also known as the IP alias, should not be confused with the base IP address, which is the IP address that corresponds to the host name of a system.
VxFS (Veritas File System)	A component of the Veritas Storage Foundation product suite that provides high performance and online management capabilities to facilitate the creation and maintenance of file systems. A file system is a collection of directories organized into a structure that enables you to locate and store files.
VxVM (Veritas Volume Manager)	A Symantec product installed on storage clients that enables management of physical disks as logical devices. It enhances data storage management by controlling space allocation, performance, data availability, device installation, and system monitoring of private and shared systems.
VVR (Veritas Volume Replicator)	A data replication tool designed to contribute to an effective disaster recovery plan.

Index

A

- agent log
 - format 186
 - location 186
- agents
 - intelligent resource monitoring 36
 - poll-based resource monitoring 36
- AMF driver 35

B

- binary message catalogs
 - about 191
 - location of 191

C

- Changing the CVM master 165
- cluster
 - Group Membership Services/Atomic Broadcast (GAB) 28
 - interconnect communication channel 25
- Cluster File System (CFS)
 - architecture 33
 - communication 33
 - overview 32
- Cluster master node
 - changing 165
- Cluster Volume Manager (CVM)
 - architecture 30
 - communication 31
 - overview 30
- commands
 - format (verify disks) 214
 - vxctl enable (scan disks) 214
- communication
 - communication stack 24
 - data flow 23
 - GAB and processes port relationship 29
 - Group Membership Services/Atomic Broadcast GAB 28
 - interconnect communication channel 25

- communication (*continued*)
 - requirements 24
- coordination point definition 54
- coordinator disks
 - DMP devices 48
 - for I/O fencing 48
- CP server
 - deployment scenarios 55
 - migration scenarios 55
- CP server database 43
- CP server user privileges 44
- CVM master
 - changing 165

D

- data corruption
 - preventing 44
- data disks
 - for I/O fencing 47
- drivers
 - tunable parameters 237

E

- engine log
 - format 186
 - location 186
- environment variables
 - MANPATH 96
- error messages
 - agent log 186
 - engine log 186
 - message catalogs 191
 - node ejection 275
 - VxVM errors related to I/O fencing 274

F

- file
 - errors in Oracle trace/log files 217
 - errors in trace/log files 217

fire drills

- about 171
- disaster recovery 171
- for global clusters 171

format command 214

G**GAB**

tunable parameters 238

GAB tunable parameters

- dynamic 240
 - Control port seed 240
 - Driver state 240
 - Gab queue limit 240
 - Halt on process death 240
 - Halt on rejoin 240
 - IOFENCE timeout 240
 - Isolate timeout 240
 - Keep on killing 240
 - Kill_ntries 240
 - Missed heartbeat halt 240
 - Partition arbitration 240
 - Quorum flag 240
 - Stable timeout 240
- static 238
 - gab_conn_wait 238
 - gab_flowctrl 238
 - gab_isolate_time 238
 - gab_kill_ntries 238
 - gab_kstat_size 238
 - gab_logbufsize 238
 - gab_msglogsize 238
 - gab_numnids 238
 - gab_numports 238

getdbac

troubleshooting script 180

I**I/O fencing**

- communication 47
- operations 47
- preventing data corruption 44
- testing and scenarios 49

intelligent resource monitoring

- disabling manually 122
- enabling manually 122

IP address

troubleshooting VIP configuration 224

K**kernel**

tunable driver parameters 237

L**LLT**

- about 25
- tunable parameters 245

LLT timer tunable parameters

setting 252

LMX

- error messages, non-critical 273
- tunable parameters 253

log files 209**logging**

- agent log 186
- engine log 186
- message tags 186

M**MANPATH environment variable 96****Master node**

changing 165

message tags

about 186

messages

- LMX error messages, non-critical 273
- node ejected 275
- VXFEN driver error messages 274

minors

appearing in LMX error messages 217

O**Oracle Disk Manager (ODM)**

overview 84

Oracle instance

definition 18

Oracle patches

applying 102

R**reservations**

description 46

S**SCSI-3 PR 46**

secure communication 78

- security 77
- server-based fencing
 - replacing coordination points
 - online cluster 158
- SF Oracle RAC
 - about 15
 - architecture 18, 21
 - communication infrastructure 23
 - error messages 271
 - high-level functionality 18
 - overview of components 22
 - tunable parameters 237
- SF Oracle RAC components
 - Cluster Volume Manager (CVM) 30
- SF Oracle RAC installation
 - pre-installation tasks
 - setting MANPATH 96
- Switching the CVM master 165

T

- troubleshooting
 - CVMVolDg 216
 - error when starting Oracle instance 222
 - File System Configured Incorrectly for ODM 226
 - getdbac 180
 - logging 186
 - Oracle log files 222
 - overview of topics 213, 216, 226
 - restoring communication after cable
 - disconnection 213
 - running scripts for analysis 180
 - scripts 180
 - SCSI reservation errors during bootup 199
 - shared disk group cannot be imported 214

V

- VCS
 - logging 186
- VCSIPC
 - errors in Oracle trace/log files 217
 - errors in trace/log files 217
 - overview 86
 - warnings in trace files 216
- Veritas Operations Manager 92
- Virtual Business Service
 - features 89
 - overview 88
 - sample configuration 90

- vxctl command 214
- VXFEN driver error messages 274
- VXFEN driver informational message 275
- VxVM
 - error messages related to I/O fencing 274
- VxVM (Volume Manager)
 - errors related to I/O fencing 274