

Oracle® Solaris Cluster Essentials

ORACLE®
SOLARIS



Tim Read

Many of the designations used by manufacturers and sellers to distinguish their products are claimed as trademarks. Where those designations appear in this book, and the publisher was aware of a trademark claim, the designations have been printed with initial capital letters or in all capitals.

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

AMD, Opteron, the AMD logo, and the AMD Opteron logo are trademarks or registered trademarks of Advanced Micro Devices. Intel and Intel Xeon are trademarks or registered trademarks of Intel Corporation. All SPARC trademarks are used under license and are trademarks or registered trademarks of SPARC International, Inc. UNIX is a registered trademark licensed through X/Open Company, Ltd.

The author and publisher have taken care in the preparation of this book, but make no expressed or implied warranty of any kind and assume no responsibility for errors or omissions. No liability is assumed for incidental or consequential damages in connection with or arising out of the use of the information or programs contained herein.

This document is provided for information purposes only and the contents hereof are subject to change without notice. This document is not warranted to be error-free, nor subject to any other warranties or conditions, whether expressed orally or implied in law, including implied warranties and conditions of merchantability or fitness for a particular purpose. We specifically disclaim any liability with respect to this document and no contractual obligations are formed either directly or indirectly by this document. This document may not be reproduced or transmitted in any form or by any means, electronic or mechanical, for any purpose, without our prior written permission.

The publisher offers excellent discounts on this book when ordered in quantity for bulk purchases or special sales, which may include electronic versions and/or custom covers and content particular to your business, training goals, marketing focus, and branding interests. For more information, please contact:

U.S. Corporate and Government Sales
(800) 382-3419
corpsales@pearsontechgroup.com

For sales outside the United States, please contact:

International Sales
international@pearsoned.com

Visit us on the Web: informit.com/ph

Library of Congress Cataloging-in-Publication Data

Read, Tim, 1963-

Oracle Solaris Cluster essentials / Tim Read.
p. cm.

Includes bibliographical references and index.

ISBN 978-0-13-248622-4 (pbk. : alk. paper)

1. Oracle Solaris Cluster. 2. Integrated software. I. Title.

QA76.76.I57R43 2011

005.5—dc22

2010027191

Copyright © 2011 Oracle and/or its affiliates. All rights reserved.

500 Oracle Parkway, Redwood Shores, CA 94065

Printed in the United States of America. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, write to:

Pearson Education, Inc.
Rights and Contracts Department
501 Boylston Street, Suite 900
Boston, MA 02116
Fax: (617) 671-3447

ISBN-13: 978-0-13-248622-4

ISBN-10: 0-13-248622-9

Text printed in the United States on recycled paper at Courier in Stoughton, Massachusetts.

First printing, September 2010



Contents

Preface	xiii
Acknowledgments	xvii
About the Author	xix
Chapter 1 Oracle Solaris Cluster: Overview	1
The Value of the Solaris Cluster Software	1
High Availability, High Performance Computing, and Fault-Tolerant Systems	3
High Availability Compared with Disaster Recovery	4
Benefits of the Solaris Cluster Framework	4
Solaris Cluster Architecture	6
Server Configurations	6
Storage Connectivity	7
Server-to-Storage Connectivity Topologies	8
How to Protect Your Data	13
Disksets, Disk Groups, and Zpools	14
Data Protection Using Storage-Based Replication	18
Public Networking	21
Configuring IPMP for Your Public Networks	22

	Solaris Cluster Public Network Monitoring Daemon	23
	Limiting Network Interface Requirements Using VLAN Tagging	24
	Network Performance: Jumbo Frames and Link Aggregation	24
	Private Networking	26
	Campus and Metro Clusters	28
	Stretching Cluster Node Separation to the Limit	28
	Performance Impact of Separating Data Centers	30
Chapter 2	Oracle Solaris Cluster: Features and Architecture	35
	Cluster Architecture	35
	Object Request Broker Subsystem	38
	High-Availability Framework	39
	Cluster Configuration Control	42
	Cluster Configuration Repository	42
	Cluster Membership	44
	Split-Brain Condition	44
	Membership	44
	Majority Voting	45
	Reconfiguration	55
	Version Manager	56
	Storage Device Fencing Subsystem	57
	Disk Fencing	57
	Disk Failfast Driver	58
	NAS Device Fencing	59
	Devices	59
	Global Devices	59
	Device ID	61
	Device Namespace	64
	Disk Path Monitoring	65
	Choosing the Correct Device to Use	66
	File Systems on a Solaris Cluster System	68
	The Cluster File System	68

SAM-QFS Software	79
Highly Available Local (Failover) File Systems	84
Cluster Public Networking	86
Logical (Virtual) Hosts	86
Global Networking Service	89
Cluster Private Interconnect	99
Private Interconnect Heartbeats	99
Private Interconnect Topology	99
Private Interconnect Traffic	100
<code>clprivnet0</code> Virtual Network Interface	101
Benefits of <code>clprivnet0</code> to Oracle RAC	102
Cluster Resilience to Private Network Failure	103
Protocols Used by the Cluster Private Network	103
TCP/IP	103
InfiniBand	104
Choosing Your Subnet Allocation for the Private Network	104
Network Ports Used by the Solaris Cluster Software	105
Configuration Guidelines	105
Data Service and Application Agents	105
Data Service Constructs	107
Resource Group Manager Daemon	107
Resource Types	108
Resources	112
Resource Groups	117
Parallel Services	123
Daemon Processes	126
Chapter 3 Combining Virtualization Technologies with Oracle Solaris Cluster Software	131
Defining a Cluster Node	132
Defining a Cluster	133
Comparison of “Black-Box” and “Fine-Grained” Control of Virtualized Entities	133

Dynamic System Domains	134
Oracle VM Server for SPARC	136
I/O Domains	138
Guest Domains	139
Failover Guest Domains	141
Oracle Solaris Zones	143
Minimal Performance Overhead	144
IP Exclusive and IP Shared Networking Options	145
Oracle Solaris Zones Root Directory	145
Oracle Solaris Zones Patch Process	146
Services That Cannot Run in Oracle Solaris Zones	146
Branded Zones (BrandZ)	146
HA Containers	147
Global-Cluster Non-Voting Node	148
Zone Clusters	150
Security Isolation	151
Application Fault Isolation	152
Resource Management	152
Dedicated Cluster Model	153
Single Point of Administration	153
Administrative Workload Reduction	154
Assigning Storage Resources	154
Zone-Cluster Architecture	154
Storage Devices	158
Networks	160
Chapter 4 Managing Your Oracle Solaris Cluster Environment	163
Installing the Oracle Solaris OS on a Cluster Node	163
Root Disk Partition Requirement for the Solaris Cluster Software	164
Planning for Upgrades	165
Securing Your Solaris Operating System	166
Operating Environment Minimization	166

Operating System Hardening	168
Securing Network Communications	169
Solaris Cluster Software Installation	169
Time Synchronization	172
Cluster Management	173
Command-Line Interface	173
The Solaris Cluster Manager Graphical User Interface	175
Solaris Cluster Wizards	177
Role-Based Access Control	178
Cluster Monitoring	180
Sun Management Center Integration	181
Solaris Cluster SNMP Management Information Base	182
Service-Level Management and Telemetry	183
Gathering Telemetry from the Solaris Cluster Software	185
Patching and Upgrading Your Cluster	189
Upgrade Methods	190
Upgrading Nodes Using Oracle Solaris Zones	196
Backing Up Your Cluster	198
Root Disk Backup	198
Backing Up Application Data on a Cluster	199
Highly Available Backup Servers	200
Creating New Resource Types	201
Application Suitability	201
Generic Data Service	203
Supporting New Applications Using the Advanced Agent Toolkit	207
Developing Resource Types by Creating a Subclass of the GDS	208
scdsbuilder GUI	212
Resource Type Registration File	216
Resource Management API	218
Data Service Development Library	218

	Useful Utilities for Building Custom Data Services	219
	Tuning and Troubleshooting	220
Chapter 5	Oracle Solaris Cluster Geographic Edition: Overview	223
	Why Have a Disaster Recovery Solution?	223
	Choosing an Appropriate Disaster Recovery Solution	224
	Benefits of a Third-Party Disaster Recovery Framework	225
	Solaris Cluster Geographic Edition Architecture	226
	Comparison of Automated and Automatic Service Migration	226
	Protecting Your Data Using Replication	227
	Storage-Based Replication: EMC Symmetrix Remote Data Facility and Hitachi Universal Replicator	228
	Host-Based Replication: StorageTek Availability Suite	230
	Application-Based Replication	232
	Protecting File Systems with Host-Based and Storage-Based Replication	233
	Connection Topologies Supported by Solaris Cluster Geographic Edition	235
	Three-Data-Center Architectures: Combining Campus and Geographic Options	237
	Using Solaris Cluster Geographic Edition with Virtualization Technologies	242
	Using Geographic Edition with Dynamic System Domains	242
	Using Geographic Edition with Oracle Solaris Zones	243
	Using Geographic Edition with Logical Domains	244
Chapter 6	Oracle Solaris Cluster Geographic Edition: Features and Architecture	247
	Software Infrastructure Required for the Geographic Edition Software	248
	Solaris Cluster Resource Groups and Resources	248
	Geographic Edition Common Agent Container Modules	249
	Event Propagation	252

CLI, GUI, and Module Implementation	252
Storage of Geographic Edition Configuration Information	252
Creating Trust between Clusters	253
Partnerships	254
Geographic Edition Heartbeat Messages	257
Heartbeat Module	258
Failure Notification	259
Protection Groups	260
Replication Components	263
StorageTek Availability Suite Software	265
EMC Symmetrix Remote Data Facility	273
Hitachi Data Systems TrueCopy and Universal Replicator	282
Oracle Data Guard for Oracle Real Application Clusters Databases	291
MySQL Replication	300
Script-Based Plug-In Module	306
Null (none) Data Replication Type	312
Protecting Oracle RAC Databases with Storage-Based Replication Products	313
Starting and Stopping Protection Groups	313
Switchover and Takeover	315
Chapter 7 Managing Your Oracle Solaris Cluster Geographic Edition Systems	321
Installing and Removing the Geographic Edition Software	321
Patching and Upgrading	324
Cluster Management	324
Command-Line Interface	325
Role-Based Access Control (RBAC)	325
Monitoring	325
Troubleshooting	327
Creating Additional Data Replication Modules	329

Chapter 8	Example Oracle Solaris Cluster Implementations	331
	Test-Driving Solaris 10 OS and Solaris Cluster Software Using Oracle VM VirtualBox Software	331
	Installing the Solaris Cluster Software to Create a Two-Node Cluster	339
	Creating a Highly Available Oracle 11g Release 1 Database	358
	Setting Up Solaris Cluster Telemetry	372
	Creating a Scalable Web Service Using Global-Cluster Non-Voting Nodes	377
	Creating an HA-Oracle Database Instance in a Zone Cluster	387
Chapter 9	Example Oracle Solaris Cluster Geographic Edition Implementations	395
	Configuring Oracle Solaris Cluster Geographic Edition	395
	Protecting a Scalable Web Service Using StorageTek Availability Suite	398
Bibliography		407
	References	407
	Additional Resources	410
Index		411



Preface

It has now been over eight years since Richard Elling and I wrote our Sun BluePrints *Designing Enterprise Solutions with Sun Cluster 3.0* [Design] to herald the release of the new Sun Cluster 3.0 software. Since then, the Solaris Cluster software has grown and matured, providing a stable and robust clustering product on which you can confidently run your business- and mission-critical services.

Now with the completion of the Oracle acquisition, the Oracle Solaris Cluster software is part of a much larger software portfolio. This software will continue to provide the best-of-breed high-availability solutions for the Oracle Solaris platform, addressing business continuity and disaster recovery for customers deploying Oracle and non-Oracle databases and applications.

Much of the additional content for this book is drawn from white papers and Sun BluePrints documents that I have authored or coauthored over the intervening years, augmented with material taken from Sun blog postings, documentation, and other Sun sources of technical information. See the complete list of sources used for this book in the Bibliography at the end of the book.

Scope

This book is written from the perspective that the Solaris 10 10/09 release is your default operating system and that you have installed, or intend to install, the Solaris Cluster 3.2 11/09 software. Because this book represents the Solaris Cluster software capabilities at the time of this writing, I have deliberately not included

many product configuration restrictions. Inevitably, some of these restrictions might be lifted as the product is tested and developed. Consequently, always check the latest product documentation, or contact Oracle technical staff, if you have any questions about a particular configuration or a product feature.

Although the book provides considerable detail on both Oracle Solaris Cluster and Oracle Solaris Cluster Geographic Edition, I have ultimately had to draw the line regarding what content to include and what content to omit. In general, I have not gone in depth into any of the Solaris Cluster agents or specific command-line options. My expectation is that this book gives you a sufficiently detailed description of the products' features and capabilities to provide a solid grounding in the subject matter. From there, you can expand your knowledge through reading the extensive product literature.

The book covers both the Oracle Solaris Cluster and Oracle Solaris Cluster Geographic Edition products. Because the Solaris Cluster software is a prerequisite for the Solaris Cluster Geographic Edition software, the book covers the Solaris Cluster software first. Chapters 1 through 4 are devoted to the Solaris Cluster software, and Chapters 5 through 7 describe the Solaris Cluster Geographic Edition software. The last two chapters, 8 and 9, provide several detailed examples of how to use the respective products.

Although every effort has been made to check the accuracy of the information in this book, it is almost inevitable that some errors have evaded the scrutiny of the book's reviewers or my corrections based on their feedback. So, I think the following quotation from one of my favorite books of all time is pertinent:

The Hitch Hiker's Guide to the Galaxy is an indispensable companion to all those who are keen to make sense of life in an infinitely complex and confusing Universe, for though it cannot hope to be useful or informative on all matters, it does at least make the reassuring claim, that where it is inaccurate it is at least *definitively* inaccurate. In cases of major discrepancy it is always reality that's got it wrong.¹

Intended Audience

This book is intended to serve as a valuable resource for three broad categories of readers. These readers span the entire lifecycle of a cluster system, from initial design through implementation and finally ongoing support.

- *Data center architects* will benefit from the material that describes the product's requirements, features, and capabilities. This book enables them to

1. Douglas Adams, *The Restaurant at the End of the Universe* (Pan paperback edition), p. 35.

design Solaris Cluster systems suited to providing the levels of availability required by the target applications.

- *System implementers* will derive the most value from the sections that describe how the Solaris Cluster resource group and resource constructs can best be used to encapsulate the target applications.
- *System administrators* will find the book a useful source of background material on how the Solaris Cluster features work. This material as well as the sections on troubleshooting will help them resolve day-to-day issues with the cluster. In addition, the sections on cluster management will provide them with guidance on maintaining a Solaris Cluster system after it has been moved into production.

The book assumes that the reader is familiar with the Oracle Solaris operating system (Oracle Solaris OS).

How to Use This Book

Readers who just want to know about the high-availability features provided by the Solaris Cluster software, and who are not interested in disaster recovery, will benefit from reading the first four chapters of the book. Chapter 1 focuses on the basic requirements for designing a Solaris Cluster system, and Chapter 2 provides a more in-depth description of the Solaris Cluster features. If you are considering implementing virtualization technologies, Chapter 3 explains how the Solaris Cluster software can coexist with these environments. Finally, Chapter 4 discusses some of the tasks that you perform to manage your Solaris Cluster system.

Readers who want to understand the capabilities of the Solaris Cluster Geographic Edition software, or who are planning a disaster recovery system and are already familiar with the Solaris Cluster software, should start with Chapter 5. In many ways, Chapter 5 mirrors Chapter 1 insofar as it provides a general background to the Solaris Cluster Geographic Edition product. More detailed technical information about the features of the Geographic Edition product is explained in Chapter 6. Finally, Chapter 7 covers some additional tasks that you might need to perform while managing a Solaris Cluster Geographic Edition installation.

If you learn best by playing with a piece of software and then reading the associated background material, then Chapter 8 and Chapter 9 are the places to start. Each chapter has detailed examples that you can follow to create your own Solaris Cluster or Solaris Cluster Geographic Edition installation. Bear in mind that these examples are not necessarily optimal configurations or what might be considered “best practice.” Instead, the examples focus on demonstrating feature capabilities to give you a platform on which you can build and experiment.

Note

Oracle Corporation acquired Sun Microsystems, Inc., early in 2010, when this book was nearing completion. Although this book mostly uses the new product names, occasional reference is also made to previous names. The following table provides a guide to the old and new product names.

Sun Product Name	Oracle Product Name
Solaris	Oracle Solaris
Solaris Cluster	Oracle Solaris Cluster
Solaris Cluster Geographic Edition	Oracle Solaris Cluster Geographic Edition
Solaris Zones or Solaris Containers	Oracle Solaris Zones or Oracle Solaris Containers
Logical Domains	Oracle VM Server for SPARC
VirtualBox	Oracle VM VirtualBox
ZFS	Oracle Solaris ZFS



4

Managing Your Oracle Solaris Cluster Environment

This chapter covers many important aspects of an Oracle Solaris Cluster system's lifecycle. It begins by outlining some of the planning you need to do prior to installing the Oracle Solaris operating system on a cluster node. Next comes information about how to secure your newly installed system against attack. The options for installing the Solaris Cluster software are then covered, followed by information about the day-to-day tasks involved in managing and monitoring the cluster. The chapter concludes by showing how to add support for new applications for which there is no third-party agent from Oracle or the application vendor.

Installing the Oracle Solaris OS on a Cluster Node

You can install the Oracle Solaris OS in these ways:

- Interactively using the `installer` program
- With Solaris JumpStart to automate the process
- With public domain or third-party automation tools such as JumpStart Enterprise Toolkit [JetWiki]

If you choose an automated option, then you can perform the installation using collections of Solaris packages or a Solaris Flash archive file that is created using the `flarcreate` command from a previous template installation (see the `flarcreate(1M)` man page) [FlarInstall].

Note

Although it is technically feasible to install a cluster node from a Flash archive that already includes preconfigured Solaris Cluster software, Oracle does not support doing so. You can, however, use an archive with the unconfigured Solaris Cluster packages added.

If your root disk is not already protected by some form of hardware RAID, you should mirror it to protect your node against disk failure. You can protect your root disk using the bundled Solaris Volume Manager software or with the Oracle Solaris ZFS file system. ZFS is available only when you use the Solaris 10 OS. If you have a software license for Veritas Volume Manager, you can also use that to protect your root disk.

Solaris Live Upgrade eases the process of upgrading your systems. Solaris Live Upgrade is available in all Solaris releases beginning with the Solaris 8 OS. However, the combination of the Solaris 10 OS, ZFS, and Solaris Live Upgrade provides the simplest and most flexible upgrade option for your systems.

When choosing hostnames for your cluster nodes, you must ensure that they comply with RFC 1123 [RFC1123]. If you intend to install Solaris Cluster Geographic Edition, you must not use an “_” (underscore) in the hostnames. Other applications might place additional constraints on the hostnames. Changing a hostname after a cluster has been installed is a complex process and should be avoided, if possible.

Although these sections on installing the Solaris Cluster software cover many important points, they do not provide all the steps you need to follow. Therefore, you must consult the latest versions of the Oracle Solaris OS [S10InstallGuide] and the Solaris Cluster software [SCInstallGuide] installation guides.

Root Disk Partition Requirement for the Solaris Cluster Software

You can use either UFS or ZFS as the root (/) file system for a Solaris Cluster node. The way you partition your root disks depends on which file system you choose and whether you also use Solaris Volume Manager.

After installation of the Solaris Cluster software, each node has a 512-megabyte `/global/.devices/node@X` file system that is mounted globally (the `X` represents the node number). To achieve this configuration, you can allow the `scinstall` program to create a `lofi`-based file system for you that works for both UFS and ZFS root (/) file systems. Alternatively, if you use UFS for the root (/) file system, you can create and mount a 512-megabyte file system on `/globaldevices` and allow the `scinstall` program to reuse it for the `/global/.devices/node@X` file system.

If you intend to use Solaris Volume Manager, then you must create the main state replica databases. These databases also require separate disk partitions, which should be 32 megabytes in size. This partition is usually placed on slice 7 of the root disk. This can pose a problem if you are using ZFS for the root disk because the standard Solaris installation uses the entire disk for the root zpool, unless you pre-partition your disk before the installation process begins. You can achieve a root disk layout where slice 7 is 32 megabytes in size, and use a ZFS root (/) file system, if you install your system using JumpStart Enterprise Toolkit (JET). However, you must not use ZFS volumes (zvols) to store the Solaris Volume Manager state replica databases, as this configuration is not supported.

Because Solaris Volume Manager relies on a state replica majority (see the section “Solaris Volume Manager’s State Replica Majority” in Chapter 1, “Oracle Solaris Cluster: Overview”) to maintain data integrity if a disk failure occurs, you should assign slices from three separate, nonshared disks. If you do not have three separate slices available, you can place all of the replicas on a single disk as shown in the following example.

Example 4.1 Creating the Solaris Volume Manager Root State Replica Databases

Use the `metadb` command to create the Solaris Volume Manager root state replica databases.

```
# metadb -afc 3 c1t0d0s7
# metadb
      flags          first blk      block count
a m  pc lu0        16             8192      /dev/dsk/c1t0d0s7
a   pc lu0        8208           8192      /dev/dsk/c1t0d0s7
a   pc lu0       16400           8192      /dev/dsk/c1t0d0s7
```

If you use Veritas Volume Manager, you ideally use it only for shared disk management, and you use either Solaris Volume Manager or the Oracle Solaris ZFS file system to mirror your root disks.

Planning for Upgrades

After you have installed your cluster, you must maintain it through its lifecycle. Doing so inevitably requires that you install both Solaris and cluster patches. You might also perform more major changes such as upgrading to a new Solaris Cluster release. To minimize any disruption, you must plan ahead.

“Patching and Upgrading Your Cluster” describes in more detail the options for ongoing cluster maintenance. If you choose ZFS for your root disk, then the upgrade procedure is fairly straightforward because Solaris Live Upgrade can create an alternate boot environment on the existing root zpool. However, if you choose UFS, then you must have the same number of slices available as your existing root disk has, which usually means using a separate disk. However, if you have only a root partition, you can plan for future upgrades by setting aside a separate partition of the same size on the root disk.

Securing Your Solaris Operating System

The security of your Solaris Cluster system must be foremost in your mind. Without security measures in place, your system could be vulnerable to denial-of-service attacks or have its data stolen, compromised, or deleted. Any of these threats could lead to your service being unavailable. Therefore, you must put barriers in place to prevent such threats.

Although several methods for auditing security violations are available, this section focuses on how to reduce the number of targets that an unauthorized user can attack. Even when your cluster is in a physically secure location, you can still take the additional measures of operating system minimization and hardening.

Note

The Solaris Cluster software does not yet support Solaris 10 Trusted Extensions.

Operating Environment Minimization

With operating system minimization, you install the fewest number of Solaris packages that are required to operate your cluster and the services it supports. The Solaris installation process allows you to choose which software groups you install. These groups are shown in Figure 4.1. Furthermore, you can customize the installation to the level of individual packages.

Clearly, choosing the software group at the center of the diagram means installing the fewest packages and thus presents the smallest “attack profile” to potential intruders. However, this installation is insufficient for the Solaris Cluster 3.2 1/09 software, which requires at least the End User Solaris Software Group (SUNWCuser) [SCInstallGuide].

One problem you face when trying to minimize your Oracle Solaris OS is that of ongoing maintenance, particularly with new or changing application require-

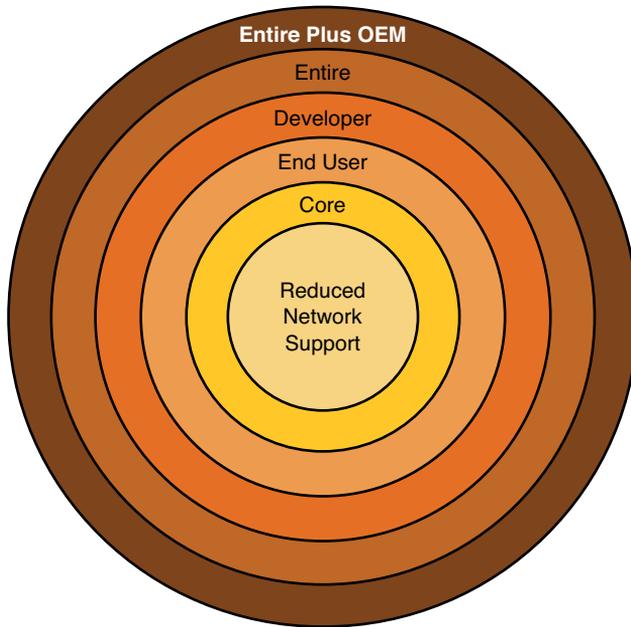


Figure 4.1 Solaris software groups for installation

ments. Though your initial system might well be served by just the `SUNWCuser` package group, you could find that a new service requires some additional packages. Not only does this mean that you need to install the missing packages, but you might also have to reapply certain patches that would have been applied to these packages had they been installed in the first place. This also assumes that the Solaris packages on which your applications are dependent are clearly documented.

For these reasons, you must take a pragmatic approach to minimization. If your systems are unlikely to change over time, and testing has confirmed that the `SUNWCuser` group is sufficient for your needs, then installing that group should pose no problems. Conversely, if the services hosted on your cluster are likely to change, then installing the `Entire` package group will better serve your needs.

Regardless of your starting point, any subsequent change to the Solaris package installation list must be performed under a strong change management control process. Such a process would require thorough testing of your services in the new environment and verifying that they can still switch over successfully. Because such testing is disruptive, any preparatory work must be performed in a test environment before final testing during a planned maintenance period.

Operating System Hardening

In contrast to operating system minimization, which reduces the number of packages installed on the system, operating system hardening disables or modifies the configuration of services that are installed on the system. An example is the disabling of the `rlogin`, `rsh`, and `rcp` commands in preference to their secure counterparts, `ssh` and `scp`.

As with Oracle Solaris OS minimization, there is clearly an almost unlimited scope for OS hardening. Unfortunately, not all changes are necessarily compatible with running the Solaris Cluster software. Furthermore, if you have made several custom changes to your system, it might prove tricky for Oracle's support services to determine whether any fault that might arise has a root cause in your hardening choices or in the Solaris Cluster software itself. Consequently, they might request that you reverse your hardening effort in order to facilitate their diagnosis. Doing so could be complicated if you do not have a well-documented change control process or an automated system for applying your hardening rules.

Fortunately, Sun has developed a supported system to help you harden your systems effectively. The package is known as the Solaris Security Toolkit (formally known as the JumpStart Architecture and Security Scripts, or JASS) [SSTAdmin-Guide]. This allows you to harden your Oracle Solaris OS in a supported and reversible manner, regardless of whether you are using the Solaris 9 OS or the Solaris 10 OS.

In addition, in the Solaris 10 OS 11/06 release comes the Secure by Default configuration [SBDBlog]. This configuration does not make the Solaris Security Toolkit obsolete, but it does subsume many of the settings that the toolkit previously changed as part of the hardening process. The latest version of the toolkit still provides fine-grained control of your hardening process and performs many other security-related tasks.

With Secure by Default configuration, you need to relax some of these settings in the Solaris 10 OS prior to installing the Solaris Cluster software. The commands you need to execute are provided in the following example [ThorstenBlog].

Example 4.2 Reversing Some of the Solaris 10 Secure by Default Settings for a Solaris Cluster Installation

Ensure that the `local_only` property of `rpcbind` is set to `false`.

```
# svcprop network/rpc/bind:default | grep local_only
```

If `local_only` is not set to `false`, run

```
# svccfg
svc:> select network/rpc/bind
svc:/network/rpc/bind> setprop config/local_only=false
svc:/network/rpc/bind> quit
# svcadm refresh network/rpc/bind:default
```

This value is needed for cluster communication between nodes.

Ensure that the `tcp_listen` property of `webconsole` is set to `true`.

```
# svcprop /system/webconsole:console | grep tcp_listen

If tcp_listen is not true, run

# svccfg
svc:> select system/webconsole
svc:/system/webconsole> setprop options/tcp_listen=true
svc:/system/webconsole> quit
# svcadm refresh svc:/system/webconsole:console
# /usr/sbin/smcwebserver restart
```

This value is needed for Solaris Cluster Manager communication.

To verify that the port is listening to `*.6789`, run

```
# netstat -an | grep 6789
```

Securing Network Communications

To secure your cluster on the network, start from the position of denying access to everything. Then proceed to allow access to ports that are specifically required by your applications and ports that are required by your management and administration procedures.

To secure network communications, you can use external firewalls. Alternatively, if you are using the Solaris 10 OS, you can use the built-in `ipfilter` capabilities. If you plan to use scalable services with shared IP addresses, you must use an external firewall. This is because the `ipfilter` software (see the `ipfilter(5)` man page) cannot detect the connection state after the shared address migrates from one cluster node to another cluster node.

In addition, if the applications on your cluster transmit sensitive data over the network to client systems, consider using the IP Security Architecture (IPsec) to protect the communication traffic (see the `IPsec(7P)` man page).

Solaris Cluster Software Installation

After you have installed and secured the Oracle Solaris OS, you can install the Solaris Cluster software. If you use Solaris JumpStart, then you can set up the JumpStart

server to install the Solaris Cluster software as part of the postinstallation scripts. This Solaris Cluster software installation process installs only the relevant packages for you and does not actually configure the software. To achieve this level of automation, you can add your own scripts to the JumpStart framework. Many installations performed by Sun Professional Services use the freely available JumpStart Enterprise Toolkit (JET). JET enables you to create individual modules to perform specific tasks such as installing or configuring software. One such module, available as part of certain Sun Professional Services engagements, can install and perform a basic configuration of the Solaris Cluster software.

The installation process has two parts: installing the software and configuring the software. The latter requires your nodes to be rebooted as part of that process. Software installation is performed through the `installer` program found on the installation media.

The `installer` program can be run in both character mode and as a full graphical user interface (GUI). To run the GUI installer program, ideally use `ssh` with X forwarding enabled. If you choose not to use `ssh`, you can use `telnet` or `rlogin` in combination with the `xhost` command to allow the cluster client node to display on the X server display you are using. If you use `telnet` or `rlogin`, you must set the `DISPLAY` environment variable to the workstation from which you are accessing the cluster.

When run from the command line, the `installer` program can be supplied with a text-based “response” file to allow automation of the Solaris Cluster software installation.

Example 4.3 Starting the Solaris Cluster Software `installer` Program

Use the `-X` flag to `ssh` to forward the X11 protocol to your local display.

```
admin_console# ssh -X root@phys-earth1
Password:
Last login: Thu Nov 12 08:19:29 2009 from admin_console.
Sun Microsystems Inc. SunOS 5.10 Generic January 2005
```

Check that `ssh` is configured to forward correctly.

```
cluster1# grep X11F /etc/ssh/sshd_config
X11Forwarding yes
```

Change to the Solaris Cluster media directory and run the installer.

```
cluster1# cd DVD-mount-point
cluster1# ls
Copyright License README Solaris_sparc Solaris_x86
```

```
cluster1# cd Solaris_sparc/
cluster1# ls -l
total 28
drwxr-xr-x  9 root    root          9 Aug  6 16:17 Product
-rwxr-xr-x  1 root    root        10892 Jan  8 2009 installer
-rw-r--r--  1 root    root          84 Jan  8 2009 release_info
cluster1# ./installer
```

Assuming you are installing from the installation media, after you have accepted the software license agreement, a wizard screen similar to Figure 4.2 is displayed.

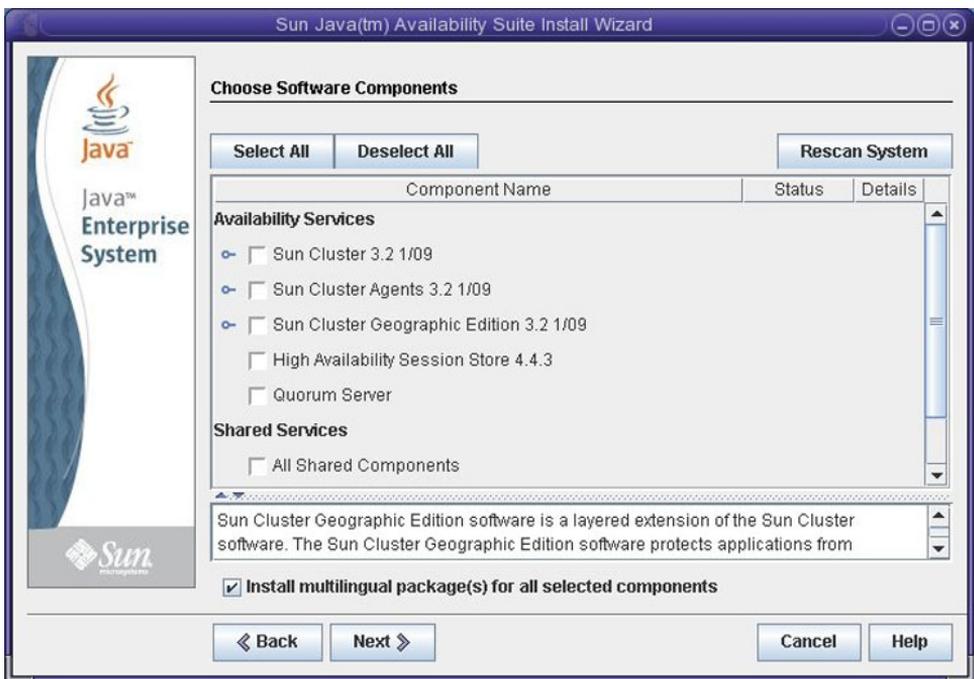


Figure 4.2 Installation screen for the Solaris Cluster 3.2 01/09 software

In the wizard, select the check boxes for the software you have chosen to install. Note that the option to install the multilingual packages is checked by default, whereas with the text `installer` program, you must explicitly choose to install these packages. Follow the wizard instructions to complete the installation. You must perform this installation on all the nodes that will form your cluster.

After the packages have been installed on all the nodes, you can proceed with the configuration of the cluster itself by using the `scinstall` command. You have several configuration options to choose from. The simplest option configures the Solaris Cluster software on just one node. The most automated option automatically configures the Solaris Cluster software on all your cluster nodes from the single command.

As each node is configured, it is rebooted by the configuration process. Therefore, do not run this program unless your nodes can experience an outage at this time.

If the configuration is performed by a certified Solaris Cluster installer, he or she will follow Oracle System Installation standards methodology to complete the installation. Doing so mostly involves checking that certain settings have been made and that the appropriate patches have been installed. If you choose to install the cluster yourself, you must purchase a cluster validation service from Sun Services to ensure that your installation complies with the necessary standards for it to be supportable.

After all the cluster nodes have been correctly configured and validated, you can begin the process of creating the resource groups and resources that encapsulate the services you choose to make highly available.

Time Synchronization

Every aspect of managing, securing, planning, and debugging a network involves determining when events happen. Time is the critical element that enables an event on one network node to be mapped to a corresponding event on another network node. In many cases, these challenges can be overcome by the enterprise deployment of the Network Time Protocol (NTP) service. You can configure this service to operate in the following modes: NTP server, NTP client, or NTP peer.

The `xntpd` daemon (see the `xntpd(1M)` man page) is bundled with the Solaris software to provide time synchronization services to all cluster nodes. The Solaris Cluster software installation creates a `/etc/inet/ntp.conf.cluster` file on each cluster node and a legacy `/etc/rc2.d/S74xntpd_cluster` script (or legacy Solaris Management Facility [SMF] service) to start `xntpd` with this file, instead of the default `/etc/inet/ntp.conf` file. The standard Solaris 10 SMF service, `/network/ntp:default`, is disabled.

The `ntp.conf.cluster` file, created by the Solaris Cluster configuration process, synchronizes the time across all cluster nodes by naming the cluster nodes as peers. You can also synchronize your cluster nodes to an external NTP server either

by naming a specific NTP server or by broadcasting or multicasting for one. To do this, you can specify the appropriate directive in the `/etc/inet/ntp.conf.cluster` file on each cluster node. Because the `S74xntpd_cluster` script calls `ntpdate` (see the `ntpdate(1M)` man page), do not restart this service while the cluster is running. If the cluster nodes are not closely synchronized with the external source, restarting the service will result in a large and sudden change in the system clocks, either backward or forward. Under some circumstances, this change can result in a system panic. Instead, follow the recommended procedure, which involves rebooting the cluster. Clearly, you should schedule such a task for the next maintenance window.

Warning

When using NTP, do not attempt to adjust the cluster time while the cluster is running. Do not adjust the time by interactively using the `date`, `rdate`, `xntpd`, or `svcadm` command or within `cron` scripts.

Cluster Management

For the most part, you can manage your cluster using either the command-line interface (CLI) or the Solaris Cluster Manager GUI. Although the capabilities provided by the GUI are a subset of what is available through the CLI, its graphical nature makes it less prone to user errors when creating or modifying complex objects. This is particularly true if a resource has many mandatory parameters or if a resource group requires several affinities to be configured.

Command-Line Interface

The Solaris Cluster CLI commands are in the `/usr/cluster/bin` directory, and all begin with the `cl` prefix. Most of these commands have both a long form and a short form. The directory also contains the CLI commands used to manage systems installed with software releases prior to the Solaris Cluster 3.2 software, such as `scrgadm`, `scconf`, and `scswitch`.

All of the commands listed in Table 4.1, excluding `clsetup`, which is menu-driven, conform to this format: `command action argument ... operand`.

All the commands return an exit code of zero on success, allowing you to create scripts to automate procedures that you might perform regularly.

Table 4.1 Solaris Cluster Commands

CLI— Long Form	CLI— Short Form	Description
<code>claccess</code>	<code>claccess</code>	Manages the Solaris Cluster access policies for nodes
<code>cldevice</code>	<code>cldev</code>	Manages the Solaris Cluster devices
<code>cldevicegroup</code>	<code>cldg</code>	Manages the Solaris Cluster device groups
<code>clinterconnect</code>	<code>clintr</code>	Manages the Solaris Cluster private interconnects
<code>clnasdevice</code>	<code>clnas</code>	Manages access to NAS devices from the Solaris Cluster software
<code>clnode</code>	<code>clnode</code>	Manages the Solaris Cluster nodes
<code>clquorum</code>	<code>clq</code>	Manages the Solaris Cluster quorum devices and properties
<code>clreslogicalhostname</code>	<code>clrslh</code>	Manages resources for the Solaris Cluster logical hostnames
<code>clresource</code>	<code>clrs</code>	Manages resources for the Solaris Cluster data services
<code>clresourcegroup</code>	<code>clrg</code>	Manages resource groups for the Solaris Cluster data services
<code>clresourcetype</code>	<code>clrt</code>	Manages resource types for the Solaris Cluster data services
<code>clressharedaddress</code>	<code>clrssa</code>	Manages the Solaris Cluster resources for shared addresses
<code>clsetup</code>	<code>clsetup</code>	Menu-driven command used to configure the Solaris Cluster software interactively
<code>clsnmphost</code>	<code>clsnmphost</code>	Administers the list of Solaris Cluster Simple Network Management Protocol (SNMP) hosts
<code>clsnmpmib</code>	<code>clsnmpmib</code>	Administers the Solaris Cluster SNMP management information bases (MIBs)
<code>clsnmpuser</code>	<code>clsnmpuser</code>	Administers the Solaris Cluster SNMP users
<code>cltelemetryattribute</code>	<code>clta</code>	Configures system resource monitoring
<code>cluster</code>	<code>cluster</code>	Manages the global configuration and status of a cluster
<code>clvsvm</code>	<code>clvsvm</code>	Configures VxVM for the Solaris Cluster software
<code>clzonecluster</code>	<code>clzc</code>	Creates and manages zone clusters

The Solaris Cluster Manager Graphical User Interface

You can run the Solaris Cluster Manager GUI using a browser such as Mozilla Firefox or Microsoft Internet Explorer. Solaris Cluster Manager requires that you have already installed

- The common agent container (CAC) package (SUNWcacaort)
- The Solaris Cluster Manager packages (SUNWscspmr, SUNWscspmu)
- The Sun Java Web Console packages (SUNWmcon, SUNWmconr, SUNWmctag)

After the Sun Java Web Console is started, it listens for secure https connections on port 6789, if you have reversed the hardening performed by the standard Solaris installation process (see Example 4.2). In addition, at least version 2.2 of the CAC service must be running. CAC listens for Java Management Extensions (JMX) calls made by the Sun Java Web Console as a result of actions you perform through the Solaris Cluster Manager GUI.

Example 4.4 Checking That the Required Components Are in Place to Run Solaris Cluster Manager

Check that the required packages are present.

```
# pkginfo SUNWcacaort
application SUNWcacaort Common Agent Container - Runtime

# pkginfo | grep SUNWscspm
application SUNWscspmr          Sun Cluster Manager (root)
application SUNWscspmu          Sun Cluster Manager (usr)

# pkginfo SUNWmcon SUNWmconr SUNWmctag
application SUNWmcon  Sun Java(TM) Web Console 3.0.2 (Core)
system              SUNWmconr Sun Java(TM) Web Console 3.0.2 (Root)
application SUNWmctag Sun Java(TM) Web Console 3.0.2 (Tags & Components)
```

Check the CAC service version.

```
# cacaoadm -V
2.2.0.1
```

Check that CAC is online.

```
# svcs /application/management/common-agent-container-1
STATE          STIME          FMRI
online         7:58:28       svc:/application/management/common-agent-container-1:default
```

Check that the Java Web Console is online.

```
# svcs svc:/system/webconsole:console
STATE      STIME      FMRI
online     Jan_19     svc:/system/webconsole:console
```

Check that the Java Web Console is listening on all networks to port 6789.

```
# netstat -a | grep 6789
*.6789          *.*          0          0 49152        0 LISTEN
```

Check the applications deployed to the Java Web Console.

```
# wcadmin list
```

Deployed web applications (application name, context name, status):

```
SunClusterManager SunClusterManager [running]
console            ROOT [running]
console            com_sun_web_ui [running]
console            console [running]
console            manager [running]
legacy             SunFlexManagerGeo [running]
zfs                zfs [running]
```

Registered jar files (application name, identifier, path):

```
SunClusterManager cacao_jars /usr/lib/cacao/lib/*.jar
SunClusterManager jdmc_jars /opt/SUNWjdmc/5.1/lib/*.jar
console            audit_jar /usr/lib/audit/Audit.jar
console            console_jars /usr/share/webconsole/lib/*.jar
console            jato_jar /usr/share/lib/jato/jato.jar
console            javahelp_jar /usr/jdk/packages/javahelp-2.0/lib/*.jar
console            shared_jars /usr/share/webconsole/private/container/shared/
lib/*.jar
```

Registered login modules (application name, service name, identifier):

```
console ConsoleLogin userlogin
console ConsoleLogin rolelogin
```

Persistent Jvm options:

```
-XX:ParallelGCThreads=4
-server
-Xmx256m
-XX:+BackgroundCompilation
-XX:+UseParallelGC
```

Shared service properties (name, value):

```
ENABLE          yes
java.options    -server -XX:+BackgroundCompilation -Xmx256m
```

After you have logged in to the Java Web Console and chosen the Solaris Cluster Manager option, a screen similar to Figure 4.3 is displayed. Using Solaris Cluster Manager, you can perform a wide range of tasks, including

- Creating, modifying, and deleting resource groups and resources
- Enabling and disabling resources
- Bringing resource groups online and offline
- Switching resource groups between nodes

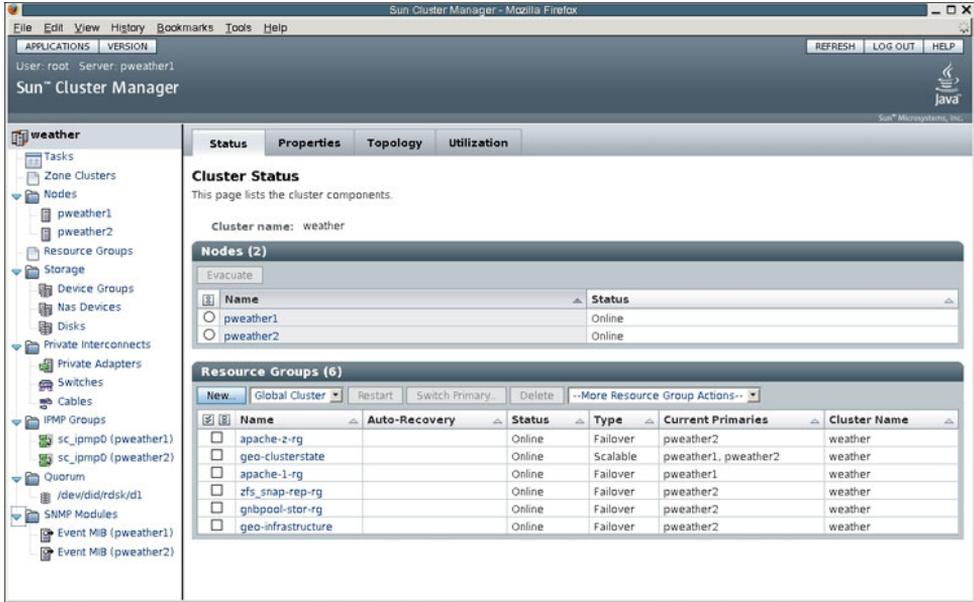


Figure 4.3 Solaris Cluster Manager browser-based GUI

Solaris Cluster Manager also highlights faults in the system with a red marker over the components that have errors.

Solaris Cluster Wizards

As with the `clsetup` menu-driven command, Solaris Cluster Manager has several wizards to help you set up some common highly available services, such as Network File System (NFS), Apache web server, and Oracle databases, as well as two more complex services: Oracle Real Application Clusters (Oracle RAC) and SAP web application server.

The wizards simplify your work by probing the system to discover the potential options that are available to you and then configuring the relevant resource groups

and resources with the appropriate dependencies and affinities. A good example of this process is the Oracle RAC wizard. It can do the following:

- Create the RAC framework resource group and resources for you, depending on your answer regarding your volume management method (hardware RAID, Solaris Volume Manager for Sun Cluster, or the VxVM cluster feature), shared QFS file system, or another supported combination.
- Create the necessary scalable storage resource groups and resources.
- Create the RAC server proxy resource group and resource based on the Oracle RAC version, the `ORACLE_HOME` variable, and the `ORACLE_SID` variables you select. Again, most of these choices are discovered for you by the wizard.
- Create the Oracle Clusterware resources that integrate with the Solaris Cluster storage resources to give the Oracle RAC database instances visibility into the availability of their underlying storage.

The wizard completes the operations for you and displays the Solaris Cluster commands it used to achieve these steps. If anything goes wrong during the creation of these objects, then the wizard will undo the changes.

Role-Based Access Control

By default, only the Solaris root user can manage the cluster. Non-root users cannot create, modify, or delete resource groups or resources, nor can they switch resource groups between the nodes. Non-root users can view the resource group and resource configuration because they are assigned the `solaris.cluster.read` rights profile as part of their standard authorizations. To delegate some of the management functions to specific system administrators who are not assigned full root access, you can use the role-based access control (RBAC) capabilities to assign them the necessary `solaris.cluster.admin` or `solaris.cluster.modify` Solaris Cluster management rights profile.

The following example shows you how to assign to the Solaris user `myadmin` all the Solaris Cluster management rights profiles, which include both `solaris.cluster.admin` and `solaris.cluster.modify`. Initially, without being assigned the `solaris.cluster.modify` rights profile, the user could not create the resource group `foo-rg`. After the user is assigned this rights profile, the user can create and delete the resource group. However, in order to manage or bring the resource group online, the user also needs to have the `solaris.cluster.admin` rights profile assigned. There is no simple way to assign roles such that individual users can control individual resource groups or resources.

have the right password. However, after being assigned the role, the user can perform the actions that the role is authorized to perform.

Example 4.6 Using an RBAC Role to Allow a User to Perform Cluster Management Tasks

Display the Cluster Management entry in the profile description database.

```
# grep "Cluster Management" /etc/security/prof_attr | head -1
Cluster Management::Sun Cluster
Management:auths=solaris.cluster.admin,solaris.cluster.modify,
.
.
.
# roleadd -P "Cluster Management" clusadm
# tail -3 /etc/user_attr
root:::auths=solaris.*,solaris.grant;profiles=Web Console Management,All;lock_after_
retries=no;min_label=admin_low;clearance=admin_high
zfssnap:::type=role;auths=solaris.smf.manage.zfs-auto-snapshot;profiles=ZFS File
System Management
clusadm:::type=role;profiles=Cluster Management
# passwd clusadm
New Password:
Re-enter new Password:
passwd: password successfully changed for clusadm
# su - myadmin
Sun Microsystems Inc. SunOS 5.10 Generic January 2005
$ su clusadm
Password:
Roles can only be assumed by authorized users
su: Sorry
$ logout
# usermod -R clusadm myadmin
# tail -3 /etc/user_attr
zfssnap:::type=role;auths=solaris.smf.manage.zfs-auto-snapshot;profiles=ZFS File
System Management
clusadm:::type=role;profiles=Cluster Management
myadmin:::type=normal;roles=clusadm
# su - myadmin
Sun Microsystems Inc. SunOS 5.10 Generic January 2005
$ su clusadm
Password:
$ clresourcegroup create bar-rg
$ clresourcegroup delete bar-rg
```

Cluster Monitoring

Most data centers use an enterprise management tool such as HP OpenView, BMC Patrol, or IBM Tivoli to monitor and manage their environment. If the data center has a large base of Sun servers, then the list might also include the Sun Management Center software. The Solaris Cluster software integrates with these tools using the Simple Network Management Protocol (SNMP), enabling administrators to monitor both the cluster nodes and the resource groups and resources they contain.

Sun Management Center Integration

The Sun Management Center software enables you to manage and monitor all the Sun servers in your data center. The Sun Management Center software is implemented as a three-tier architecture: console, server, and agent. The agent is installed on the system you choose to monitor. It sends alerts to the server when certain events occur or when thresholds are exceeded. Each agent handles the monitoring of a specific function. For example, the Solaris Cluster agent monitors the cluster infrastructure and the resource groups and resources configured on your system. The Sun Management Center console aggregates the alarms and alerts pertaining to a specific node, marking the node with a symbol to denote the most serious condition that it is encountering. If there are no problems with the node, then no symbol is present. Figure 4.4 shows an example of a serious error with the cluster transports, but just a warning for the scalable resource group because it is offline.

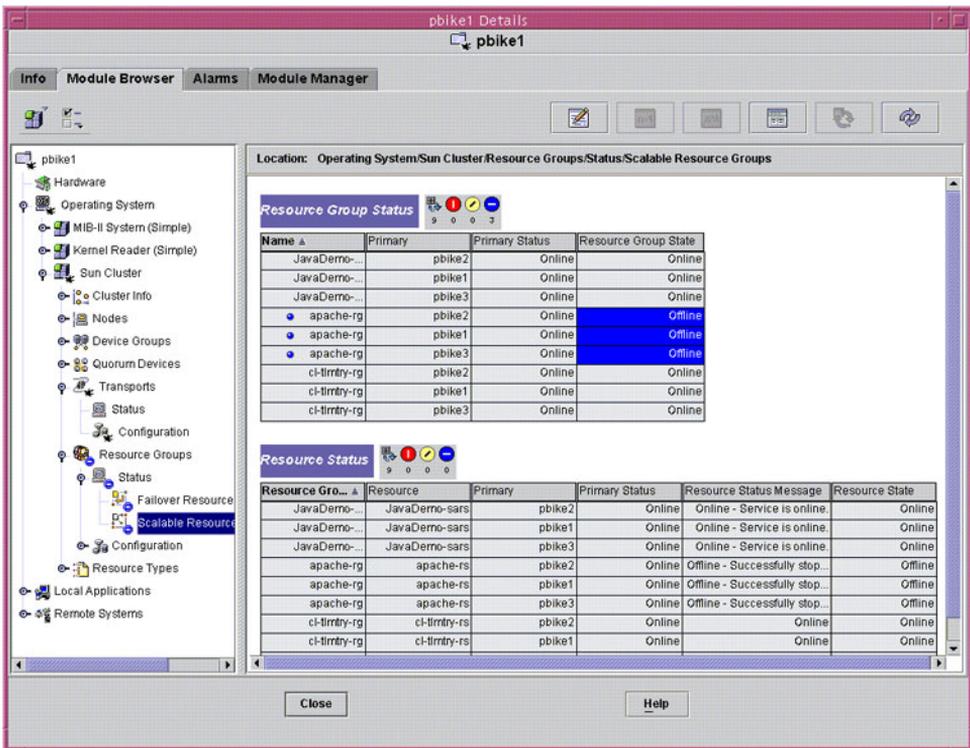


Figure 4.4 Sun Management Center GUI

Solaris Cluster SNMP Management Information Base

The Solaris Cluster software provides two Simple Network Management Protocol (SNMP) management information bases (MIBs): one for cluster events, the other for cluster objects and their status.

SNMP MIB for Cluster Events

The Solaris Cluster SNMP event MIB can be enabled on any cluster node without relying on additional SNMP agent infrastructure. Because all cluster nodes receive the same cluster events, you need to enable the SNMP event MIB on only one cluster node.

The cluster events are graded in severity from low to high, as follows:

- INFO
- WARNING
- ERROR
- CRITICAL
- FATAL

When the SNMP event MIB is enabled on a node, the cluster event SNMP interface can do the following:

- Store up to 50 of the most recent WARNING or higher-severity cluster events in the MIB table.
- Send SNMP traps to one or more hosts when WARNING or higher-severity cluster events occur.

By default, the SNMP event module uses port number 11161, and traps use port number 11162. These port numbers can be changed by using the `cacoadm` command, as shown in the following example.

The Solaris Cluster SNMP event MIB Object ID (OID) is 1.3.6.1.4.1.42.2.80.

Example 4.7 Changing the SNMP Adapter Trap Port

Use the `cacoadm` command to determine and change the SNMP adapter trap port.

```
phys-grass2# cacoadm get-param snmp-adaptor-trap-port
snmp-adaptor-trap-port=11162
phys-grass2# cacoadm stop
phys-grass2# cacoadm set-param snmp-adaptor-trap-port=12345
phys-grass2# cacoadm start
phys-grass2# cacoadm get-param snmp-adaptor-trap-port
snmp-adaptor-trap-port=12345
```

SNMP Interface for Cluster Objects

The SNMP interface for cluster objects such as nodes, resource groups, and quorum devices, along with their corresponding status, is provided through the Solaris Cluster software to the Sun Management Center integration module. The integration module is installed by default on all cluster nodes as part of the `SUNWscsa` and `SUNWscsam` packages.

To enable the integration module on a cluster node, the Sun Management Center agent infrastructure must be installed and running on that node. After the Sun Management Center agent starts on the node, the integration module can be loaded and enabled through the Sun Management Center console. The Sun Management Center server and console can be installed on one of the cluster nodes, but they are usually installed on another system.

Similar to the SNMP event MIB, the integration module needs to be enabled on only one cluster node to receive cluster SNMP MIB data and traps. The integration module MIB OID is `1.3.6.1.4.1.42.2.80.1.1.1`.

The Sun Management Center agent running on a cluster node sends SNMP traps to the Sun Management Center server whenever the status of an object changes, such as a resource group going online or offline, or becoming managed or unmanaged. The SNMP traps also can be sent directly from the Sun Management Center agent to any other hosts by adding the secondary trap destination through the Sun Management Center `es-trapdest` command. The following example shows how to add `myHost` to the hosts that receive Sun Management Center SNMP traps from the agent on port 11162.

Example 4.8 Adding a Secondary Trap Destination to the Sun Management Center Agent

Add `myHost` as a secondary trap destination using the `es-trapdest` command.

```
# /opt/SUNWsymon/sbin/es-trapdest -c agent -a myHost 11162
```

The agent has to be restarted for the setting to take effect.

Service-Level Management and Telemetry

When you consolidate multiple services onto a Solaris Cluster installation, you must ensure that your service levels are met even when several services reside on the same cluster node. The Oracle Solaris OS has many features, such as resource controls and scheduler options, to help you achieve these goals. These resource

allocations can be defined in the projects database stored locally in `/etc/project` or held in the name service maps.

The Solaris Cluster software can bind both resource groups and resources to projects using the `RG_project_name` and `Resource_project_name` properties, respectively. The following example shows how to create a processor pool (containing four CPUs) that uses the fair share scheduler (FSS). The processor pool is then associated with the user project that limits shared memory usage to 8 gigabytes. The FSS can be enabled by using the `dispadmin -d FSS` command.

Example 4.9 Binding a Resource Group to a Project Associated with a Processor Pool

Determine the number of processors the system has using the `psrinfo` command.

Define a four-CPU-processor set called `oracle_pset` in a temporary file, and then use the file as input to the `poolcfg` command.

```
# psrinfo | wc -l
24
# cat /tmp/create_oracle_pool.txt
create pset oracle_pset ( uint pset.min = 1 ; uint pset.max = 4 )
create pool oracle_pool
associate pool oracle_pool ( pset oracle_pset )
modify pool oracle_pool ( string pool.scheduler = "FSS" )

# poolcfg -f /tmp/create_oracle_pool.txt
```

Instantiate the configuration using the `pooladm` command.

```
# pooladm -c
# pooladm

system default
  string system.comment
  int system.version 1
  boolean system.bind-default true
  string system.poold.objectives wt-load

pool pool_default
  int pool.sys_id 0
  boolean pool.active true
  boolean pool.default true
  string pool.scheduler FSS
  int pool.importance 1
  string pool.comment
  pset pset_default

pool oracle_pool
  int pool.sys_id 2
  boolean pool.active true
  boolean pool.default false
  string pool.scheduler FSS
  int pool.importance 1
  string pool.comment
  pset oracle_pset
```

```

pset oracle_pset
    int      pset.sys_id 1
    boolean  pset.default false
    uint     pset.min 1
    uint     pset.max 4
    string   pset.units population
    uint     pset.load 17
    uint     pset.size 4
    string   pset.comment

    cpu

    int      cpu.sys_id 1
    string   cpu.comment
    string   cpu.status on-line

    cpu

    int      cpu.sys_id 0
    string   cpu.comment
    string   cpu.status on-line

    cpu

    int      cpu.sys_id 3
    string   cpu.comment
    string   cpu.status on-line

    cpu

    int      cpu.sys_id 2
    string   cpu.comment
    string   cpu.status on-line

pset pset_default
    int      pset.sys_id -1
    boolean  pset.default true
.
.
.

```

Use the `projadd` command to make `oracle_pool` the project pool for user `oracle`.

```

# projadd -p 4242 -K "project.max-shm-memory=(privileged,8GB,deny)" \
> -K project.pool=oracle_pool user.oracle
# su - oracle
Sun Microsystems Inc. SunOS 5.10 Generic January 2005
$ id -p
uid=424242(oracle) gid=424242(oinstall) projid=4242(user.oracle)
$ exit
# clresourcegroup create -p RG_project_name=user.oracle oracle-rg

```

Similarly, using the `clzonecluster` command (see the `clzonecluster(1M)` man page), you can bind zone clusters to pools, dedicate or limit the number of CPUs allocated to them, and limit the physical, swap, or locked memory they can use.

Gathering Telemetry from the Solaris Cluster Software

The Solaris Cluster service-level management feature enables you to configure the Solaris Cluster software to gather telemetry data from your cluster. Using this

feature, you can collect statistics on CPU, memory, swap, and network utilization of the cluster node as well as on resource groups and system components such as disks and network adapters. By monitoring system resource usage through the Solaris Cluster software, you can collect data that reflects how a service using specific system resources is performing. You can also discover resource bottlenecks, overloads, and even underutilized hardware resources. Based on this data, you can assign applications to nodes that have the necessary resources and choose which node each application should fail over to.

This feature must be set up using the `clsetup` command. The telemetry data is stored in its own Java DB database held on a failover or global file system that you must provide for its use. After the setup is complete, you can enable the telemetry on the resource groups, choose the attributes to monitor, and set thresholds. Figure 4.5 and Figure 4.6 show the type of output you can receive from using this feature.

Figure 4.5 shows that an alarm has been generated because disk `d4` has exceeded the threshold set for it.

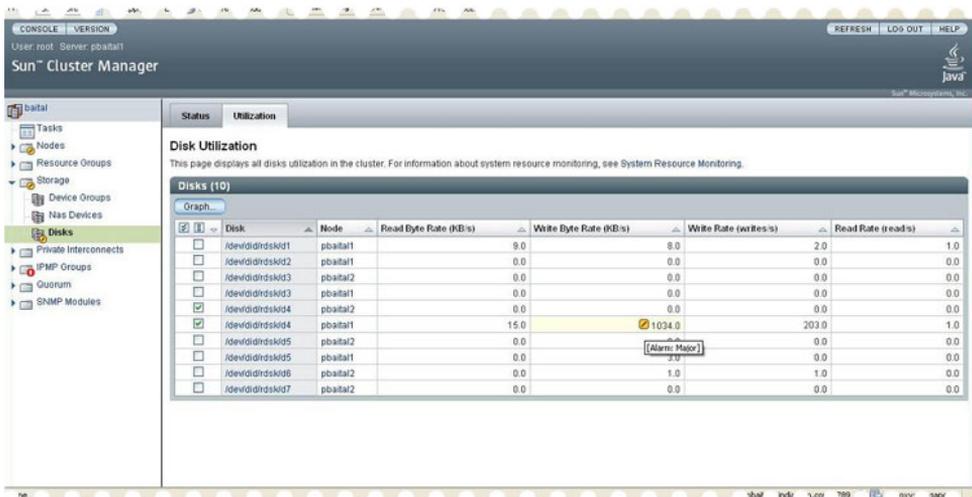


Figure 4.5 Alarm showing that the write I/O rate to disk `d4` has exceeded the threshold set

Figure 4.6 shows the utilization of the public network adapters `bge0` and `bge1` on cluster node `pbaital1`.

The telemetry uses the `RG_slm_type` resource group property, which can be set to one of two values: `automated` or `manual`. The default value for the `RG_slm_type` property is `manual`. Unless the `RG_slm_type` property value is explicitly set to

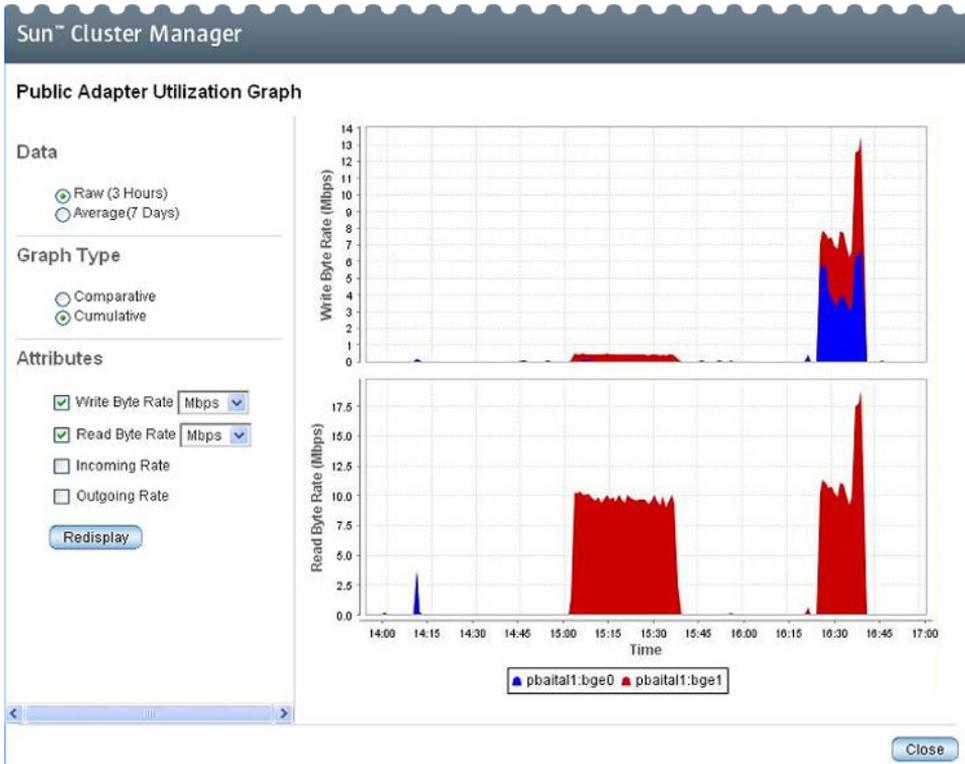


Figure 4.6 Public network adapter utilization telemetry gathered using the service-level management feature

automated when a resource group is created, telemetry is not enabled for the resource group. If the resource group `RG_slm_type` property is changed, resource utilization monitoring begins only after the resource group is restarted.

When a resource group has the `RG_slm_type` property set to `automated`, the Resource Group Manager (RGM) internally generates a Solaris project to track the system resource utilization for all processes encapsulated by the resource of the resource group. This tracking happens regardless of whether the `RG_project_name` and `Resource_project_name` properties are set. The telemetry can track only the system resource utilization: CPU usage, resident set size (RSS), and swap usage for resource groups that have the `RG_slm_type` property set to `automated`. Telemetry for other objects is gathered at the node, zone, disk, or network interface level, as appropriate.

See Example 8.9 in Chapter 8, “Example Oracle Solaris Cluster Implementations,” for more information about how to set up, configure, and use the Solaris Cluster telemetry.

Using the Solaris Cluster Manager browser interface simplifies the process of configuring thresholds and viewing the telemetry monitoring data.

The following example shows the generated project name in the `RG_SLM_projectname` property. However, unlike other resource group properties, you cannot set this property manually. Furthermore, if `RG_slm_type` is set to `automated`, the `RG_project_name` and `Resource_project_name` properties will be ignored. Conversely, when `RG_slm_type` is set to `manual`, the processes of the resource group’s resource will be bound to the projects named in the `RG_project_name` and `Resource_project_name` properties. However, the RGM will not track the system resources they use.

Example 4.10 The Effect of Setting the `RG_slm_type` Property to `automated`

Use the `clresourcegroup show -v apache-1-rg` command to show the property settings for the `apache-1-rg` resource group.

```
# clresourcegroup show -v apache-1-rg
=== Resource Groups and Resources ===

Resource Group:                               apache-1-rg
RG_description:                               <NULL>
RG_mode:                                       Failover
RG_state:                                      Managed
RG_project_name:                              default
RG_affinities:                                <NULL>
RG_SLM_type:                                  manual
Auto_start_on_new_cluster:                   False
Failback:                                     False
Nodelist:                                     phys-winter1 phys-winter2
Maximum primaries:                            1
Desired primaries:                            1
RG_dependencies:                              <NULL>
Implicit_network_dependencies:                 True
Global_resources_used:                        <All>
Pingpong_interval:                            3600
Pathprefix:                                   <NULL>
RG_System:                                    False
Suspend_automatic_recovery:                  False

--- Resources for Group apache-1-rg ---
.
.
.
```

Use the `clresourcegroup` command to set the `RG_SLM_type` property to `automated`.

```
# clresourcegroup set -p RG_SLM_type=automated apache-1-rg
# clresourcegroup show -v apache-1-rg

=== Resource Groups and Resources ===

Resource Group:                apache-1-rg
RG_description:                <NULL>
RG_mode:                       Failover
RG_state:                      Managed
RG_project_name:              default
RG_affinities:                <NULL>
RG_SLM_type:                  automated
RG_SLM_projectname:          SCSLM_apache_1_rg
RG_SLM_pset_type:            default
RG_SLM_CPU_SHARES:          1
RG_SLM_PSET_MIN:             0
Auto_start_on_new_cluster:    False
Failback:                     False
Nodelist:                     phys-winter1 phys-winter2
Maximum primaries:           1
Desired primaries:           1
RG_dependencies:              <NULL>
Implicit_network_dependencies: True
Global_resources_used:        <All>
Pingpong_interval:           3600
Pathprefix:                   <NULL>
RG_System:                    False
Suspend_automatic_recovery:   False

--- Resources for Group apache-1-rg ---
.
```

Patching and Upgrading Your Cluster

The approach you take to patching and upgrading your cluster will vary depending on your data center standards. Some data centers apply only the bare minimum of changes, such as critical security patches, operating on the principle of “if it ain’t broke, don’t fix it.” Clearly, if a system is stable, this is a reasonable approach, especially if availability is your highest priority. Other data centers schedule regular planned outages to perform patching, often choosing to lag a little behind the most recent patches to avoid encountering bad patches. Still others take the “latest and greatest” approach, given that the patches are meant to fix all the key, known issues. Whichever strategy you use, minimizing the overall outage will clearly be your main priority.

You can also upgrade the Solaris Cluster software by patching it with the patches that are built into a later release of the Solaris Cluster software. However, you must follow all the instructions in the README file contained within each patch.

Upgrade Methods

Three main methods are available for patching or upgrading your operating system and Solaris Cluster software: standard upgrade, rolling upgrade, and dual-partition upgrade. Solaris Live Upgrade provides yet more options when used with these methods. Each method results in differing lengths of outage on your cluster, which further depends on the number and type of zones you have configured (see the section “Oracle Solaris Zones” in Chapter 3, “Combining Virtualization Technologies with Oracle Solaris Cluster Software”). The following sections describe these methods.

Before you begin making changes to the system, consider the following checklist:

- Has the upgrade procedure been tested on a test cluster before being implemented on the production system? If not, is management aware of the potential risk of untested procedures?
- Have you verified that the new configuration is supported by all your services? Some applications place greater restrictions on the operating system or other software versions they support.
- Have the proposed changes been agreed upon and recorded in your change control system?
- Do you have a full backup of the cluster prior to making any changes?
- Have you tested that you can switch over all services between cluster nodes prior to your changes? This testing isolates any post-upgrade problems to the upgrade and patches.
- Is there an agreed outage period of a set length?
- Is there a procedure to reverse any changes if the upgrade does not go according to plan? Have you determined when you will need to start backing out any changes if that happens?
- Have you tested that you can switch over all services between cluster nodes after the changes have been made? This testing ensures that if a failure occurs later, your system will continue to provide a set of highly available services.

Standard Upgrade

With standard upgrade, you schedule a complete cluster outage to perform the patch or upgrade process. This method might be considered for applying patches, such as kernel patches, that need to be installed in single-user mode. Major operating system upgrades also fall into this category. You cannot run a Solaris Cluster configuration with nodes running different major versions of the Oracle Solaris OS, for example, Solaris 9 OS and Solaris 10 OS.

A standard upgrade is by far the simplest method but necessitates the longest cluster outage. During this period, your services are unavailable. Consequently, this method is the least desirable from an availability standpoint.

A standard upgrade can be combined with Solaris Live Upgrade. This approach significantly reduces the cluster outage by no longer requiring the entire cluster to be shut down while you are upgrading the software. You first create an alternative boot environment with the upgraded software on each machine (see “Solaris Live Upgrade” for details). When you are ready to activate the new software, halt the entire cluster, and then boot the entire cluster using the alternative boot environment with the new software on each machine.

Rolling Upgrade

With rolling upgrade, the nodes in your cluster are upgraded one at a time. Any service hosted on a node to be upgraded is manually switched over to another cluster node prior to the target node being shut down for patching. After the target node has been patched, it is brought back into the cluster, and the process is repeated until all nodes have been upgraded. After all of the nodes have been upgraded, you execute the `scversions -c` command to commit the cluster to use the new software. The cluster software does not enable cluster-wide protocols until the commit happens. Prior to the commit, all nodes, including nodes with new software, continue to use the old cluster-wide protocols. There is no rolling downgrade after the commit happens. The outage incurred by a service is therefore just twice the time it takes to switch over from one cluster node to another. A switchover of a simple application can take about 30 seconds, though the time varies significantly depending on the application.

A rolling upgrade changes only the Solaris Cluster software or the Oracle Solaris OS. Rolling upgrades can be used as long as the major cluster release and major operating system remain the same, for example, the Solaris Cluster 3.2 software or the Solaris 10 OS. You cannot perform rolling upgrades between major cluster or major OS releases, for example, the Sun Cluster 3.1 software to the Solaris Cluster 3.2 software or the Solaris 9 OS to the Solaris 10 OS.

The rolling upgrade process updates the software on the local storage of a node when that node is not in cluster mode. The highly available (HA) container places the software on shared storage. Because the shared storage is never out of service, the rolling upgrade process has no opportunity to change the software residing on shared storage. Therefore, a rolling upgrade cannot be used to upgrade software on HA containers.

A rolling upgrade can also be combined with Solaris Live Upgrade. This combination does not reduce the service outage. A rolling upgrade without Solaris Live Upgrade reduces the service outage by performing software upgrades while the

machines are not in the cluster. During this time period, the ability of the cluster to survive machine failures is reduced. The use of Solaris Live Upgrade can significantly reduce the time period of increased vulnerability. Solaris Live Upgrade also reduces the period of time in which nodes in the cluster have different software versions. You first create an alternative boot environment with the upgraded software on each machine (see “Solaris Live Upgrade”). You then serially halt each machine and boot that machine from the alternative boot environment with the upgraded software. After you have booted all nodes from the alternative boot environment with the upgraded software, you execute the `scversions -c` command to commit the cluster to use the new software.

Dual-Partition Upgrade

A dual-partition upgrade offers the greatest flexibility in terms of what changes you can make to the system. The change can be a patch, an update release, or a major release. Changes in multiple areas can be applied simultaneously. A dual-partition upgrade can even be used when all the changes are in areas other than the Solaris Cluster software, including the following:

- The Oracle Solaris OS
- A Solaris Cluster software release
- Third-party volume managers and file systems
- Application software (as long as the applications are not installed on shared storage)

A dual-partition upgrade divides your cluster into two parts, each of which is updated separately and sequentially. To invoke a dual-partition upgrade, you must use the `scinstall` command from the target media of the release to which you are upgrading. This command then prompts you with the options shown in Example 4.11. Choose option 3 for a dual-partition upgrade. You must divide the cluster machines into two sets of nodes, which are called *partitions*. You must partition the cluster such that each partition has for each storage device at least one node that can access that storage device. During the upgrade process, each partition will host the cluster applications, and the cluster applications must be able to access storage. If you are unsure of how to partition the cluster, you can choose an option whereby the system presents you with the various ways to group your cluster nodes for the purposes of the upgrade. After you have chosen your partitioning, you next select the option to initiate the actual upgrade.

Example 4.11 Main Menu of the `scinstall` Command Located on the Solaris Cluster Release Media

The main menu of the `scinstall` command provides you with the following options:

```
*** Main Menu ***
```

```
Please select from one of the following (*) options:
```

- 1) Create a new cluster or add a cluster node
- 2) Configure a cluster to be JumpStarted from this install server
- * 3) Manage a dual-partition upgrade
- * 4) Upgrade this cluster node
- * 5) Print release information for this cluster node

- * ?) Help with menu options
- * q) Quit

```
Option: 3
```

The system software makes changes to the quorum votes, migrates your services to the partition that is still running while the first partition is upgraded, and halts the nodes of the first partition. You boot the nodes of the first partition into non-cluster mode and upgrade the software. After you have completed the upgrade on the first partition, you choose to continue the dual-partition upgrade. The system software reboots the nodes of the first partition into cluster mode as a new cluster. Unlike what happens in an ordinary boot, the nodes of the first partition proceed through the boot process and stop just prior to the point where the cluster would activate resources used by applications (such as mounting file systems and plumbing IP addresses) and start applications. At this point, most of the boot process has completed. The system software halts the nodes of the second partition and waits for them to actually halt. Next, the system software imports volumes, mounts file systems, plumbs IP addresses, and starts applications. In effect, the system performs a switchover of all applications and their resources from the second partition to the new partition. Next, reboot each node of the second partition and upgrade the software. When done with that step, you choose to continue the upgrade. The system software reboots the nodes of the second partition into cluster mode. After all nodes are back in the cluster, the system software restores the quorum configuration and completes the dual-partition upgrade.

If the dual-partition upgrade fails, you boot all nodes in non-cluster mode and execute the `scinstall -u recover` command on each node to restore the initial copy of the Cluster Configuration Repository (CCR). The system is already completely shut down at this point, so a minimal downtime upgrade is not possible. You do not restart the dual-partition upgrade. Instead, you perform a standard upgrade.

The service outage when you perform a dual-partition upgrade on a system with failover application is approximately the same as for a rolling upgrade. Because a time comes when neither partition can host an instance of a scalable application, there is a service outage for the scalable application, and the service outage is approximately the same as that of a failover application. For the simplest scalable applications, the service outage can be about 30 seconds, though the time varies significantly based on the application.

You cannot upgrade any software or data on the shared storage because the services remain running throughout the upgrade process, except for the point when the services are switched over from the partition running the old software to the partition running the new software. Do not make any configuration changes during the dual-partition upgrade. Make configuration changes either before or after the dual-partition upgrade.

A dual-partition upgrade can also be combined with Solaris Live Upgrade. This combination does not reduce the service outage. A dual-partition upgrade without Solaris Live Upgrade reduces the service outage by performing software upgrades while the machines are not in the cluster. During this time period the ability of the cluster to survive machine failures is reduced. The use of Solaris Live Upgrade can significantly reduce the time period of increased vulnerability. You first create an alternative boot environment with the upgraded software on each machine (see “Solaris Live Upgrade”). You then initiate the actual dual-partition upgrade. The system software automatically manages all of the work of the dual-partition upgrade. There no longer is any need to boot into non-cluster mode, so the number of reboots is reduced by half. The net result is that the time required for upgrades is dramatically reduced.

Choosing Which Upgrade Method to Use

The choice of an upgrade method is governed by three factors: the simplicity of the upgrade process, the compatibility of the current and target software environments, and the level of availability your services must achieve during the upgrade process. Thus, if you require simplicity above all else, then the standard upgrade is the best method. If your current and target software releases do not span a major release, and you want to maintain service availability even for scalable services, then you can use a rolling upgrade. Finally, if changes to your cluster environment involve software spanning major releases and you need to minimize your service outage, then you must choose a dual-partition upgrade.

Solaris Live Upgrade

Solaris Live Upgrade is not a stand-alone technology for upgrading a Solaris Cluster system. Solaris Live Upgrade must be used as part of either a standard upgrade, a rolling upgrade, or a dual-partition upgrade. Solaris Live Upgrade

makes it possible to create an alternate boot environment with the upgraded software while the machine continues to run unaffected using the current boot environment. See “Standard Upgrade,” “Rolling Upgrade,” or “Dual-Partition Upgrade” for details on when to bring the alternate boot environment with the updated software into operational status. The boot environment with the old software is not damaged when you boot the alternate boot environment with the new software. If you encounter a problem with the new software, you can recover by halting all nodes of the cluster and then booting the entire cluster from the boot environment with the old software and configuration information. Solaris Live Upgrade is available for all the operating system releases supported by the Solaris Cluster 3.2 software.

If you are using the Solaris 10 OS with a ZFS root (/) file system, then the effort involved in creating an alternate boot environment is virtually zero because the required new file systems can be created from the existing root zpool. However, if you are using a Solaris Volume Manager or a Veritas Volume Manager mirrored root disk, then the process is significantly more complicated because you must find additional partitions to match the existing root partition allocation. Consequently, using the Solaris 10 OS and ZFS for the root disk has significant advantages. However, do not upgrade the software version of the zpools in your cluster until you are satisfied that your new environment is working correctly. If you do upgrade the zpool version, you will not be able to fail back to your old boot environment, as the old environment might not be able to mount the zpools that already have a higher version number. Such zpools cannot be downgraded.

After you have updated the Oracle Solaris OS on your new boot environment, you can proceed with updating the Solaris Cluster software. When using Solaris Live Upgrade, you must follow the documentation on how to handle specific steps relating to the volume managers, Solaris Volume Manager and Veritas Volume Manager, and which version of the Java Runtime Environment (JRE) is required.

The process for updating the Solaris Cluster software on the alternate boot environment is outlined in the following steps:

1. Create an `installer` program state file to drive the upgrade of the shared components on the alternate boot environment, but without actually performing the actions.
2. Run the `installer` program again using this state file but targeted on the alternate boot environment.
3. Update the Solaris Cluster framework by running the `scinstall` command from the new Solaris Cluster media, again targeting the alternate boot environment.
4. Update the Solaris Cluster agents by running the `scinstall` command from the new Solaris Cluster media, again targeting the alternate boot environment.

Upgrading Nodes Using Oracle Solaris Zones

When upgrading a cluster node that uses the Solaris 10 OS, you need to consider the impact the upgrade process has on any Oracle Solaris Zones (also known as Oracle Solaris Containers) that are configured on the system. These containers represent additional, virtual Solaris instances that must also be upgraded.

To maintain service availability during an upgrade, a resource group must have a node list that contains at least one physical cluster node that will still be available to host the resource group while the rolling upgrade or dual-partition upgrade is performed. This is feasible both for resource groups that use global-cluster non-voting nodes and for resource groups existing in zone clusters. However, HA containers are, by their very nature, the sole place in which the service resides. Therefore, taking the cluster node into single-user mode to patch it results in the HA container, and consequently the service, being unavailable for the duration of the update process [ZonePatch].

One way to minimize HA container service outage is to use Solaris Live Upgrade. The procedure is not particularly practical if you have tens or hundreds of HA containers because it involves many steps.

The following procedure is meant for environments where each HA container has its `zonpath`, on its own Oracle Solaris ZFS file system, in its own `zpool`, for example, `ora-pool/ora1` for `ora1`. Both `ora-pool` and `ora-pool/ora1` must have appropriate mount points set.

Solaris Live Upgrade requires full control over the `zpool`s when the Oracle Solaris OS reboots in order to perform the necessary renaming of the ZFS file systems and the changes to the zones' configuration files. Therefore, the ZFS file systems must not be under Solaris Cluster software control during the shutdown and boot process. This can be achieved by creating an alternate boot environment where each zone and its associated `zonpath` are under cluster control. Then, prior to shutting down the system, you can offline and unmanage the resource groups controlling the HA containers and import each `zpool` that holds a `zonpath`. Finally, you can copy all the affected configuration files in the CCR and the `zpool.cache` file to the alternate boot environment.

In general, this procedure requires you to perform the following steps:

1. Alter the HA container's resource group node list to just the current node.
2. Perform a Solaris Live Upgrade on the node.
3. Disable the resources of the HA container's resource group, which takes the container offline and unmanages the resource groups.
4. Import the `zpool` holding the zone's `zonpath`.

5. Copy the `zpool.cache` file and corresponding CCR files for the resource group to the alternate boot environment.
6. Shut down and reboot the node.
7. After Solaris Live Upgrade has performed all of the necessary reconfigurations, manage and bring the resource group online.
8. Detach the HA container from its other nodes, and then use Solaris Live Upgrade and reboot the other nodes in turn.
9. Attach the `zonepath` of the HA container with the `-F` (force) option to bring the container's state back to `installed`.
10. Reset the HA container's resource group node list.
11. Check that switchover still works and finish by switching the HA container back to its primary machine.

Although this procedure requires a shorter outage, other factors might make it impractical. First, the HA container cannot be failed over while the upgrade is under way because the node list has been restricted to the single node. Additionally, the service-level agreements (SLAs) for other HA containers that are deployed on the system might limit when the outage can occur because any upgrade will affect them, too. If there isn't a single time period when groups of containers can be updated, then the restriction will prevent any of the changes from being made.

The final option for upgrading HA containers is to use the `update on attach` feature in the Solaris 10 10/08 release. This option can be used only to update the packages in a container and does not allow them to be downgraded. If you use this option, you must ensure that the process does not introduce any packages that are older than the current version the Solaris Cluster software requires. If it does, you must update them again using the `installer` program.

An outline of the steps for using the `update on attach` feature follows:

1. Evacuate (or switch over) all the containers that will not be patched from the target node (`nodeT`) using `clresourcegroup switch ...`, leaving just the containers to be patched.
2. On node `nodeT`, use `zoneadm -z ... detach` to detach each remaining container from `nodeT` so that the operating system will not attempt to patch them.
3. On node `nodeT`, patch `nodeT` and reboot back into the cluster.
4. Use `clresource disable ...` to disable the HA container resources that were switched over to the alternative node (`nodeA`), and switch their resource groups back to `nodeT`. This switches over the `zonepaths` to `nodeT`.

5. On node `nodeA`, detach the container from `nodeA` so that the operating system will not attempt to patch the container during the patch process or reboot the container after the upgrade.
6. On node `nodeT`, use `zoneadm -z ... attach -u` to reattach and patch the HA container on `nodeT`.
7. Test the containers by booting and halting them manually.
8. Place the containers back under cluster control by reenabling the HA container resource.
9. Repeat the upgrade process with `nodeA`.

Ensure that HA containers are prevented from failing between the nodes while this procedure is being followed.

Backing Up Your Cluster

Although your data might be protected by hardware RAID or host-based mirroring software and even possibly replicated to another site for disaster recovery purposes, you must have a consistent, usable backup of the data on your cluster. The requirement is twofold, involving backup of the root disk and backup of the application data. Both have their own specific challenges.

Root Disk Backup

Your root disk contains the Oracle Solaris OS with numerous configuration files that the system requires to perform its tasks. Not all of these files are static. Many of the log files you need to retain for auditing and debugging purposes are highly dynamic. Therefore, you must achieve a consistent backup of your system so that you can restore your system successfully, if the need arises.

When using UFS for the root disk, only two methods are available for achieving a guaranteed consistent backup of the root file system partitions:

- Boot the system into single-user mode.
- Use both `lockfs` and `fssnap` while the system is at its normal run level.

Obviously, booting a node into single-user mode requires that you switch over all the services hosted on this node. Not only does this result in service outages, but it also means that the application might have to share the resources on its new host node, which might degrade its performance somewhat. The `lockfs/fssnap` option

seems better. However, they can result in the system pausing while the data is flushed from the buffer cache and a consistent view is reached. If this pause is too long, it might have an adverse effect on the cluster framework. Furthermore, any real-time process prevents `fssnap` from being able to lock the file system. Thus, with a Solaris Cluster installation, you must temporarily suspend the `xntpd` daemon. However, other processes, such as the Oracle 10g Real Application Clusters or Oracle 11g Real Application Clusters frameworks might make this approach unworkable.

After you have performed the backup, you can delete the snapshot and move on to the next partition on the root disk.

Example 4.12 Using `lockfs` and `fssnap` to Create a Consistent Root (/) File System Snapshot

Stop the `xntpd` daemon before locking the root (/) file system with the `lockfs` command.

```
# /etc/rc2.d/S74xntpd.cluster stop
# lockfs -f
```

Take a snapshot of the root (/) file system using the `fssnap` command before restarting the `xntpd` daemon.

```
# time fssnap -o backing-store=/spare_disk /
/dev/fssnap/0

real    0m19.370s
user    0m0.003s
sys     0m0.454s
# /etc/rc2.d/S74xntpd.cluster start
Perform backup...
# fssnap -d /dev/fssnap/0
Deleted snapshot 0.
```

For an Oracle Solaris ZFS file system, the situation is much more straightforward. By issuing a `zfs snapshot` command, you can create a consistent view of a file system that you can back up and restore with confidence. Using the `-r` flag allows you to create these snapshots recursively for all file systems below a certain mount point, further simplifying the process.

Backing Up Application Data on a Cluster

The first challenge with backing up application data when a service resides on a cluster is determining which cluster node the service is currently running on. If a

failure has recently occurred, then the service might not be running on its primary node. If you are running Oracle RAC, the database is probably running on multiple nodes simultaneously. In addition, the data might be stored on raw disk or in Oracle's Automatic Storage Management (ASM), rather than in a file system. Consequently, any backup process must be capable of communicating with the node that currently hosts the application, rather than depending on the application being on a particular node, and potentially using application-specific backup procedures or software.

Although `fssnap` can be used in certain circumstances to achieve a consistent view of the root (`/`) file system partitions for backup, do not use it with failover UFS file systems. The pause in file system activity while the snapshot is being taken might result in the service fault probe detecting a fault and causing a service failover. Furthermore, `fssnap` cannot be used with global file systems (see the section "The Cluster File System" in Chapter 2, "Oracle Solaris Cluster: Features and Architecture") because `fssnap` must be run on the UFS mount point directly and works closely with the in-memory data structures of UFS. This means that the PxFs client and server (master) must interpret the `fssnap ioctl` system calls, but this capability is not currently present in PxFs.

Once more, the Oracle Solaris ZFS snapshot feature enables you to obtain a consistent view of the application data and so is a simpler option if there are no specific tools for consistently backing up the application data.

Many backup products are available from Oracle and from third-party sources. Many have application-specific integration features, for example, the ability to integrate with Oracle's RMAN backup function. Most products can back up data stored in any file system (UFS, ZFS, QFS, VxFS) that you might have configured in your cluster.

Highly Available Backup Servers

It's obviously very important to perform regular, secure backups of your critical systems. This, in turn, means that the systems performing the backup must be sufficiently highly available. Otherwise, they might not be able to complete a backup within the time window available. Although there is little you can do to make an individual tape drive more available, you can have tape libraries housing multiple tape drives. Then the problem of availability rests with the system that controls the backups.

A backup (master) server contains the backup configuration information: catalogs of previous backups, schedules for subsequent backups, and target nodes to be backed up. Just like any other service, this collection of data files and the programs that access it can be made highly available. Thus, a highly available service can be achieved by placing the configuration files on a highly available file system,

hosted by one or more Solaris Cluster nodes, and encapsulating the backup server program in a suitable resource in a resource group.

The most common data center backup configuration uses SAN-attached tape libraries with multiple tape drives. You configure the master server to manage the backup by communicating with the client software installed on each target cluster node to be backed up. Instead of defining an entire physical server as a target, you use the logical host of the individual services that require their data to be backed up. The master server then contacts the appropriate physical node when the time comes to back up the data. If you need to back up the individual nodes, then you define the backup so that it covers only the file systems that constitute the root (/) file system. When the time comes to perform the backup, the master server directs the client to stream the necessary dataset to one or more tapes in the library.

Solaris Cluster agents are available for both the StorageTek Enterprise Backup software and Veritas NetBackup. If a Solaris Cluster agent is not available for your backup software, you can easily create one, as described in the next section.

Creating New Resource Types

As described in the section “Data Service and Application Agents” in Chapter 2, “Oracle Solaris Cluster: Features and Architecture,” Oracle has a substantial list of supported agents that cover most of the applications in your data center. These application agents are maintained by Oracle and are extensively tested on each new release of both the Solaris Cluster software and the application itself. Even so, inevitably you will have an application that is not part of the existing agent portfolio.

Application Suitability

Before creating a resource type for your application, you must determine whether the application meets the criteria for being made highly available. The following list highlights the main points you must consider. For a complete list see “Analyzing the Application for Suitability” in [SCDevGuide].

- Is your application crash-tolerant? This is important because in a highly available environment your application must be able to recover its data consistency without requiring manual intervention. If the application did require such intervention, then most of the benefits of a high-availability framework would be lost.
- Does your application rely on the physical node name of the machine, such as that resulting from calls to `uname`, `gethostbyname`, or equivalent interfaces?

If so, then when the application moves to another cluster node, the dependency on the physical hostname will probably cause the application to fail. There is a work-around to this problem, which is to interpose the `libscho.st.so.1` library. However, this work-around can sometimes raise support issues with application vendors.

- Can your application run on a multihomed system, that is, one with several public networks? Your application must be able to handle situations where IP addresses are configured and unconfigured from network adapters as services move around the cluster. This has consequences for the way your application binds to the network.
- Does your application use hard-coded path names for the location of its data? If so, then symbolic links might not be sufficient to ensure that the data is stored in a location that is compatible with using a failover or global file system. If the application renames a data file, it can break the symbolic links.

After you have determined that your application is suitable for being made highly available, you have several ways to achieve the necessary integration:

- You can use the Generic Data Service (GDS) directly and just supply the required parameters. Although you cannot define any new extension properties for the resource type you create, it is by far the simplest option.
- You can create a subclass of the GDS to create a completely new resource type. This option enables you to define one or more extension properties for your new resource type. This option is relatively simple and yet provides considerable flexibility.
- You can extend the GDS using the Advanced Agent Toolkit. Although this option does not create a new resource type, it does enable you to define one or more extension properties. This option is also relatively simple and provides considerable flexibility.
- You can use the GUI `scdsbuilder` tool and customize the resulting shell script or C source using the Resource Management API (RMAPI) and the Data Service Development Library (DSDL) APIs. If significant customization work is needed, this option might result in an increased maintenance burden.
- You can use the RMAPI or DSDL APIs directly to develop your resource type from scratch. This option trades the development and maintenance costs for ultimate flexibility and performance.

Each option is discussed in more detail in the following sections.

Generic Data Service

The Generic Data Service (GDS) is provided with the Solaris Cluster software. The `SUNW.gds` agent is packaged in the `SUNWscgds` package, which is installed as standard by the Solaris Cluster software installer program. The `SUNW.gds` agent is considered the preferred way to create both failover and scalable resources. The GDS is supported by Oracle, but you must support the script that you provide for the `Start_command`, `Stop_command`, `Probe_command`, and `Validate_command` methods.

By default, the `SUNW.gds` resource type is not registered, so you must register it before attempting to create a resource of that type. The commands in the following example show how to determine if the resource type is registered and then how to register it, if it is not already present.

Example 4.13 Registering the `SUNW.gds` Resource Type

Use the `clresourcetype` command to determine whether the `SUNW.gds` resource type needs to be registered.

```
# clresourcetype list | grep SUNW.gds
# clresourcetype register SUNW.gds
# clresourcetype list | grep SUNW.gds
SUNW.gds:6
```

In addition to the standard resource properties, the GDS agent has four properties to enable you to integrate your application: `Start_command`, `Stop_command`, `Probe_command`, and `Validate_command`. These properties are described in “Integrating Your Application-Specific Logic.” By using the GDS as the basis for your application, you automatically benefit from all the patches and feature upgrades that the GDS receives.

Example 4.14 shows how you can use the GDS to make the X11 program `xeyes` highly available. You begin by creating a `Start_command` program. In this example, a script calls the full path name of the program with a parameter that is passed to the shell script. This script must exist on all the cluster nodes on which the application is intended to run.

Next, having checked that the `SUNW.gds` resource type is registered, you create the resource group. In this example, you allow the resource group’s node list to default to all the cluster nodes.

Next, you create a resource to represent your program. In the example, the `Start_command` property is specified by the script you wrote (and which must exist on all nodes). The display parameter to use is also specified. Because this

program does not listen on any network ports, you set the `network_aware` property to `false`. This means that the probe mechanism used will be the continued existence of the `xeyes` process that the `Start_command` program leaves running in the background. By default, any resource you create is enabled so that when the resource group is brought online, the resource is automatically started. To change the default, you can specify the `-d` argument to the `clresource create` command.

The last two steps instruct the RGM that it needs to control or manage the `xeyes-rg` resource group and then to bring that resource group online. The action of bringing the resource group online starts the resource because it was created in an enabled state.

Assuming you have allowed remote X11 clients to display on your X server using `xhost` and you have specified the correct X display to use (substitute a value suited to your environment for `myhost:1.0`), then the `xeyes` program will appear on your display. You can switch the resource group between nodes and the RGM will kill the `xeyes` process and restart it on the new node, `phys-summer2`, as the example shows.

Example 4.14 Creating a Simple, Highly Available `xeyes` Service

List the script that will be used to start the `xeyes` command.

```
# cat /tmp/start_xeyes
#!/bin/ksh
/usr/openwin/demo/xeyes -display $1 &
exit 0
```

Check that the `SUNW.gds` resource type is registered, and then create the resource group and resource that will control the `xeyes` service.

```
# clresourcetype list | grep SUNW.gds
SUNW.gds:6
# clresourcegroup create xeyes-rg
# clresource create -t SUNW.gds \
> -p start_command="/tmp/start_xeyes myhost:1.0" \
> -p network_aware=false \
> -g xeyes-rg xeyes-rs
```

Use the `clresourcegroup` command to bring the `xeyes-rg` resource group online.

```
# clresourcegroup manage xeyes-rg
# clresourcegroup online xeyes-rg
# clresourcegroup status xeyes-rg
```

```
=== Cluster Resource Groups ===
```

Group Name	Node Name	Suspended	Status
xeyes-rg	phys-summer1	No	Online
	phys-summer2	No	Offline

```
# clresourcegroup switch -n phys-summer2 xeyes-rg
# clresourcegroup status xeyes-rg

=== Cluster Resource Groups ===

Group Name      Node Name      Suspended      Status
-----
xeyes-rg        phys-summer1   No              Offline
                phys-summer2   No              Online
```

To demonstrate how the GDS handles application failure, quit the `xeyes` program from your X display. You will notice that the RGM restarts the application almost instantaneously. The messages in `/var/adm/messages` (see Example 4.15) indicate that the RGM recognized the failure and restarted the service.

After the fault probe determines that the service is online, indicated by `Service is online` in `/var/adm/messages`, kill the process again. The resource has two properties that determine how many times it is restarted by the RGM within a certain time period. These properties are `Retry_count` and `Retry_interval` (see Example 4.16). After the specified number of failures, the built-in logic of the GDS determines that the current node is unhealthy and releases the service so that it can be started on another node. If the service also experiences problems on this node, then the RGM will not fail the service back to its original node unless the time period, in seconds, as defined by the resource group's `Pingpong_interval` property, has passed. Instead, the GDS attempts to keep the service running on the remaining node. This behavior is governed by another property called `Failover_mode`.

The purpose of the `Pingpong_interval` property is to prevent a service that fails to start from endlessly looping, resulting in the service migrating back and forth between cluster nodes. In a test environment, you might need to reset the value of `Pingpong_interval` to a lower value. Doing so enables you to restart your service once you have corrected any problems you encountered.

Example 4.15 Sample RGM Messages

The `/var/adm/messages` file contains information on the state changes of the resource groups and resources in the cluster.

```
Nov 23 04:00:23 phys-summer2 Cluster.RGM.global.rgmd: [ID 529407 daemon.notice]
resource group xeyes-rg state on node phys-summer2 change to RG_ONLINE
Nov 23 04:01:23 phys-summer2 Cluster.RGM.global.rgmd: [ID 922363 daemon.notice]
resource xeyes-rs status msg on node phys-summer2 change to <Service is online.>
Nov 23 04:01:25 phys-summer2 Cluster.PMF.pmfmd: [ID 887656 daemon.notice] Process:
tag="xeyes-rg,xeyes-rs,0.svc", cmd="/bin/sh -c /tmp/start_xeyes myhost:1.0", Failed
to stay up.
```

```

Nov 23 04:01:25 phys-summer2 Cluster.RGM.global.rgmd: [ID 784560 daemon.notice]
resource xeyes-rs status on node phys-summer2 change to R_FM_FAULTED
Nov 23 04:01:25 phys-summer2 Cluster.RGM.global.rgmd: [ID 922363 daemon.notice]
resource xeyes-rs status msg on node phys-summer2 change to <Service daemon not
running.>
Nov 23 04:01:25 phys-summer2 SC[,SUNW.gds:6,xeyes-rg,xeyes-rs,gds_probe]: [ID 423137
daemon.error] A resource restart attempt on resource xeyes-rs in resource group
xeyes-rg has been blocked because the number of restarts within the past Retry_
interval (370 seconds) would exceed Retry_count (2)
Nov 23 04:01:25 phys-summer2 SC[,SUNW.gds:6,xeyes-rg,xeyes-rs,gds_probe]: [ID 874133
daemon.notice] Issuing a failover request because the application exited.
Nov 23 04:01:25 phys-summer2 Cluster.RGM.global.rgmd: [ID 494478 daemon.notice]
resource xeyes-rs in resource group xeyes-rg has requested failover of the resource
group on phys-summer2.
Nov 23 04:01:25 phys-summer2 Cluster.RGM.global.rgmd: [ID 423291 daemon.error] RGM
isn't failing resource group <xeyes-rg> off of node <phys-summer2>, because there are
no other current or potential masters
Nov 23 04:01:25 phys-summer2 Cluster.RGM.global.rgmd: [ID 702911 daemon.error]
Resource <xeyes-rs> of Resource Group <xeyes-rg> failed pingpong check on node <phys-
summer1>. The resource group will not be mastered by that node.
Nov 23 04:01:25 phys-summer2 SC[,SUNW.gds:6,xeyes-rg,xeyes-rs,gds_probe]: [ID 969827
daemon.error] Failover attempt has failed.
Nov 23 04:01:25 phys-summer2 SC[,SUNW.gds:6,xeyes-rg,xeyes-rs,gds_probe]: [ID 670283
daemon.notice] Issuing a resource restart request because the application exited.

```

Example 4.16 Retry, Failover Mode, and Ping-pong Interval Properties

Use the `clresource` command to determine the property values of the `xeyes-rs` resource.

```

# clresource show \
> -p retry_count, retry_interval, failover_mode xeyes-rs

=== Resources ===

Resource:                                xeyes-rs

--- Standard and extension properties ---

Retry_interval:                          370
  Class:                                  standard
  Description:                             Time in which monitor attempts to
restart a failed resource Retry_count times.
  Type:                                    int

Retry_count:                              2
  Class:                                  standard
  Description:                             Indicates the number of times a
monitor restarts the resource if it fails.
  Type:                                    int

Failover_mode:                            SOFT
  Class:                                  standard
  Description:                             Modifies recovery actions taken
when the resource fails.
  Type:                                    enum

```

```
# clresourcegroup show -p pingpong_interval xeyes-rg
=== Resource Groups and Resources ===
Resource Group:                               xeyes-rg
Pingpong_interval:                             3600
```

In the preceding example, the display variable property can be changed only by stopping the resource and modifying the `Start_command` property. Although of little importance here, because the `xeyes` program must be restarted to change the target X server on which it displays, it does make a difference in instances where a variable can be changed while a service is running. Examples include changing debugging levels to use and changing directories for log files.

To create a resource type that has new extension properties that can be changed when you need to change them, you need to either write your resource type from scratch or create a subclass of the GDS, as described in a later section.

Supporting New Applications Using the Advanced Agent Toolkit

Many application agents in the current Solaris Cluster software release are derived from the Advanced Agent Toolkit methodology [AdvGDSTlkit]: HA-PostgreSQL, HA-MySQL, and HA containers, to name three. All three use the `SUNW.gds` agent as their basis. However, in its raw form, the `SUNW.gds` agent has some limitations.

The rationale behind the toolkit is that all new application agents have many common requirements:

- They might require one or more extension properties.
- They must provide debugging information.
- They might need to disable the process-monitoring facility (`pmfadm`) for applications that leave no obvious child processes to monitor.
- They must supply a `Start_command` script, as a minimum, and possibly `Stop_command`, `Probe_command`, and `Validate_command` scripts.

The toolkit also simplifies much of the work needed to handle Oracle Solaris Zones and SMF. Thus, providing this extended framework enables your developers to focus on the application-specific integration work rather than on debugging the framework itself. After the work is complete, the new resource type is registered using a registration script.

Developing Resource Types by Creating a Subclass of the GDS

The advantage of creating a subclass of the GDS, rather than writing a new resource type from scratch, is that the new resource type inherits all the best practices that are already part of the standard GDS code. In addition, creating a subclass of the GDS enables you to create your own resource type extension properties while retaining the same level of flexibility as if you had started from scratch. Finally, your new resource type, which is a subclass of the GDS, has a distinct name, enabling you to easily distinguish resources of the new resource type. If you instead used the Advanced Agent Toolkit or the `SUNW.gds` agent, then you would have to determine what the resource is by examining the extension properties or reviewing the code. This step would be necessary because the resource type would be set to `SUNW.gds`, rather than `MYCORP.appsvr`, for example.

You create a subclass of the GDS by creating a resource type registration (RTR) file where the `RT_basedir` parameter is set to the directory containing binaries used by the standard GDS methods: `Start`, `Stop`, `Validate`, and so on. You then extend the RTR file by defining your own resource type extension properties. Finally, you set the method parameters in the RTR file to point to your scripts that override the standard GDS behavior.

Several existing Sun resource types are implemented this way, including the HA-Logical Domain agent (`SUNW.ldom`), which was covered in the section “Failover Guest Domains” in Chapter 3, “Combining Virtualization Technologies with Oracle Solaris Cluster Software.”

The RTR file for the `SUNW.ldom` resource type is shown in Example 4.17. In this RTR file, the `RT_basedir` parameter is set to the standard directory for the GDS package, that is, `/opt/SUNWscgds/bin`. Of the standard methods, only `Init`, `Boot`, and `Validate` have been overridden using programs that are located in the `../../SUNWscxvm/bin` directory. Unlike a standard GDS resource type, the `Start_command`, `Stop_command`, `Probe_command`, and `Validate_command` properties are assigned fixed values and cannot be changed. This is indicated by the `Tunable = NONE` settings. Furthermore, each command, apart from `validate_command`, is called with a consistent set of arguments, namely, `-R %RS_NAME -T %RT_NAME -G %RG_NAME`. The `%variable` construct is similar to the `$variable` syntax found in shell scripts. It means that when a resource of this type is instantiated, use the names you assigned it as arguments. For example, if you wrote a resource type called `FOO.bar` and then created a resource group called `whizz-rg` containing a resource called `bang-rs` of this type, the argument passed would be `-R bang-rs -T FOO.bar -G whizz-rg`. With these arguments, you can then make calls to the RMAPI or DSDL APIs to retrieve or set properties.

In contrast to the `Start_command`, `Stop_command`, and `Probe_command` properties, the `Validate_command` property does not use this construct. Instead, the

RGM passes the validate command all the properties listed for the resource type on the command line. Then the validate command parses this list and determines whether the configuration is valid.

Example 4.17 RTR File for the SUNW.1dom Resource Type

The following text shows some of the key parts of the RTR file for the SUNW.1dom resource type:

```

.
.
.
RESOURCE_TYPE = "ldom";
VENDOR_ID = SUNW;
RT_DESCRIPTION = "Sun Cluster HA for xVM Server SPARC Guest Domains";

RT_version = "1";
API_version = 10;

RT_basedir=/opt/SUNWscgds/bin;

Init          =      ../../SUNWscxvm/bin/init_xvm;
Boot         =      ../../SUNWscxvm/bin/boot_xvm;

Start        =      gds_svc_start;
Stop         =      gds_svc_stop;

Validate     =      ../../SUNWscxvm/bin/validate_xvm;
Update      =      gds_update;

Monitor_start =      gds_monitor_start;
Monitor_stop  =      gds_monitor_stop;
Monitor_check =      gds_monitor_check;

Init_nodes = RG_PRIMARYES;
Failover = FALSE;

# The paramtable is a list of bracketed resource property declarations
# that come after the resource-type declarations
# The property-name declaration must be the first attribute
# after the open curly of a paramtable entry
#
# The following are the system defined properties. Each of the system defined
# properties have a default value set for each of the attributes. Look at
# man rt_reg(4) for a detailed explanation.
#
{
    PROPERTY = Start_timeout;
    MIN = 60;
    DEFAULT = 300;
}
{
    PROPERTY = Stop_timeout;
    MIN = 60;
    DEFAULT = 300;
}
.
.
.

```

```

# This is an optional property. Any value provided will be used as
# the absolute path to a command to invoke to validate the application.
# If no value is provided, The validation will be skipped.
#
{
    PROPERTY = Validate_command;
    EXTENSION;
    STRING;
    DEFAULT = "";
    TUNABLE = NONE;
    DESCRIPTION = "Command to validate the application";
}

# This property must be specified, since this is the only mechanism
# that indicates how to start the application. Since a value must
# be provided, there is no default. The value must be an absolute path.
{
    PROPERTY = Start_command;
    EXTENSION;
    STRINGARRAY;
    DEFAULT = "/opt/SUNWscxvm/bin/control_xvm start -R %RS_NAME -T %RT_NAME -G
    %RG_NAME";
    TUNABLE = NONE;
    DESCRIPTION = "Command to start application";
}

# This is an optional property. Any value provided will be used as
# the absolute path to a command to invoke to stop the application.
# If no value is provided, signals will be used to stop the application.
#
# It is assumed that Stop_command will not return until the
# application has been stopped.
{
    PROPERTY = Stop_command;
    EXTENSION;
    STRING;
    DEFAULT = "/opt/SUNWscxvm/bin/control_xvm stop -R %RS_NAME -T %RT_NAME -G
    %RG_NAME";
    TUNABLE = NONE;
    DESCRIPTION = "Command to stop application";
}

# This is an optional property. Any value provided will be used as
# the absolute path to a command to invoke to probe the application.
# If no value is provided, the "simple_probe" will be used to probe
# the application.
#
{
    PROPERTY = Probe_command;
    EXTENSION;
    STRING;
    DEFAULT = "/opt/SUNWscxvm/bin/control_xvm probe -R %RS_NAME -G %RG_NAME -T
    %RT_NAME";
    TUNABLE = NONE;
    DESCRIPTION = "Command to probe application";
}

# This is an optional property. It determines whether the application
# uses network to communicate with its clients.
#
{
    PROPERTY = Network_aware;
    EXTENSION;
    BOOLEAN;
    DEFAULT = FALSE;
}

```

```

        TUNABLE = AT_CREATION;
        DESCRIPTION = "Determines whether the application uses network";
    }

# This is an optional property, which determines the signal sent to the
# application for being stopped.
#
{
    PROPERTY = Stop_signal;
    EXTENSION;
    INT;
    MIN = 1;
    MAX = 37;
    DEFAULT = 15;
    TUNABLE = WHEN_DISABLED;
    DESCRIPTION = "The signal sent to the application for being stopped";
}

# This is an optional property, which determines whether to failover when
# retry_count is exceeded during retry_interval.
#
{
    PROPERTY = Failover_enabled;
    EXTENSION;
    BOOLEAN;
    DEFAULT = TRUE;
    TUNABLE = WHEN_DISABLED;
    DESCRIPTION = "Determines whether to failover when retry_count is exceeded
during retry_interval";
}

# This is an optional property that specifies the log level GDS events.
#
{
    PROPERTY = Log_level;
    EXTENSION;
    ENUM { NONE, INFO, ERR };
    DEFAULT = "INFO";
    TUNABLE = ANYTIME;
    DESCRIPTION = "Determines the log level for event based traces";
}

{
    Property = Debug_level;
    Extension;
    Per_node;
    Int;
    Min = 0;
    Max = 2;
    Default = 0;
    Tunable = ANYTIME;
    Description = "Debug level";
}

{
    Property = Domain_name;
    Extension;
    String;
    Minlength = 1;
    Tunable = WHEN_DISABLED;
    Description = "LDoms Guest Domain name";
}

```

```

{
    Property = Migration_type;
    Extension;
    Enum { NORMAL, MIGRATE };
    Default = "MIGRATE";
    Tunable = ANYTIME;
    Description = "Type of guest domain migration to be performed";
}

{
    PROPERTY = Plugin_probe;
    EXTENSION;
    STRING;
    DEFAULT = "";
    TUNABLE = ANYTIME;
    DESCRIPTION = "Script or command to check the guest domain";
}

{
    PROPERTY = Password_file;
    EXTENSION;
    STRING;
    DEFAULT = "";
    TUNABLE = WHEN_DISABLED;
    DESCRIPTION = "The complete path to the file containing the target host
password";
}

```

scdsbuilder GUI

To customize an agent beyond what is permitted by the GDS, you can use the Agent Builder command, `scdsbuilder` (see the `scdsbuilder(1HA)` man page). This command has three code generation options, and the resulting files are wrapped in a Solaris package that you can install on your cluster nodes:

- DSDL code (see the section “Data Service Development Library”).
- ksh code, including all the necessary `scha_control` commands (see the section “Resource Management API”). With the ksh code, you are creating your own resource type.
- A ksh registration script for a GDS agent. Here, the code generates the appropriate `clresource create` command.

You can customize the resulting code to your specific needs. However, with the ksh registration script for the GDS agent, the scope for modification is limited. The example in Figure 4.7 shows the use of the third option.

The `scdsbuilder` command starts the Solaris Cluster Agent Builder GUI, as shown in Figure 4.7. In this example, data has already been specified for each field available to the user. A short code of SUNW is specified for the vendor name, and `tstgds` is specified for the application name. This data is then used to generate

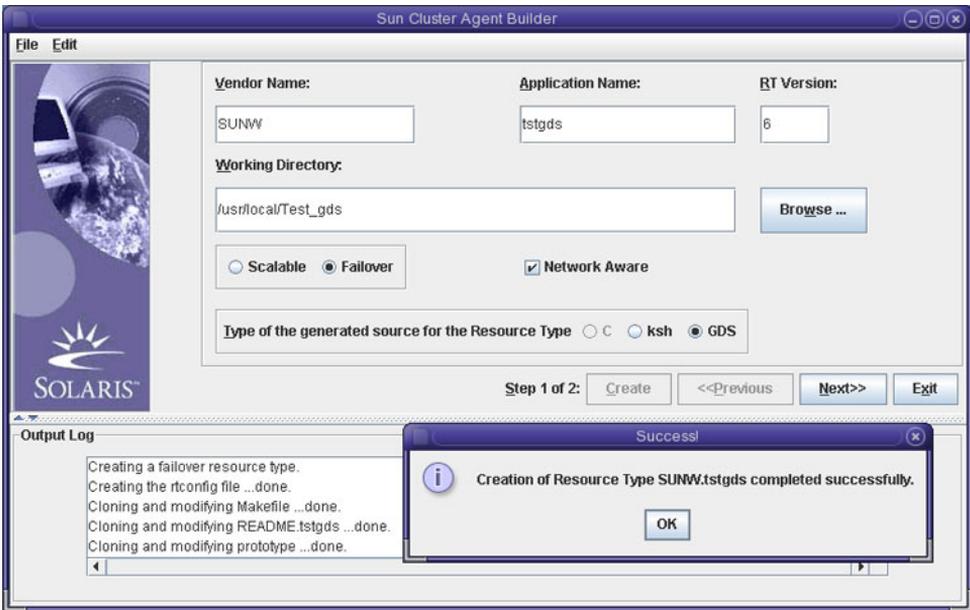


Figure 4.7 Using the scdsbuilder GUI to create a new resource type

both the name of the package that Agent Builder creates for you and the name of the resource type that you will subsequently use.

The information you provide in the other fields is used as follows:

- The RT version enables you to specify a version number for this resource type. You can identify which version of the agent you are running when it is placed into production.
- The working directory is used by Agent Builder as a working area in which it can create your package and write other associated, intermediate files.
- Your target application determines whether you select the scalable or failover option. If a particular instance of an application can run on multiple nodes at once without corrupting any of its data files, then you can select the scalable option. A good example of such an application is a web server. For all other applications, such as databases and file services, select the failover option.
- The Network Aware check box is used to determine whether any resource created using this resource type needs to have the `port_list` property set. The `port_list` property is then used by the GDS service to provide a simple probe mechanism.

- The source type option determines whether the resulting code uses the C programming language, ksh, or the GDS (see the section “SUNW.gds” in Chapter 2, “Oracle Solaris Cluster: Features and Architecture”) to create the data service. To use the C option, you must have a C compiler installed on your system.

After you have entered the data and clicked on the Next button, you are presented with the screen shown in Figure 4.8.

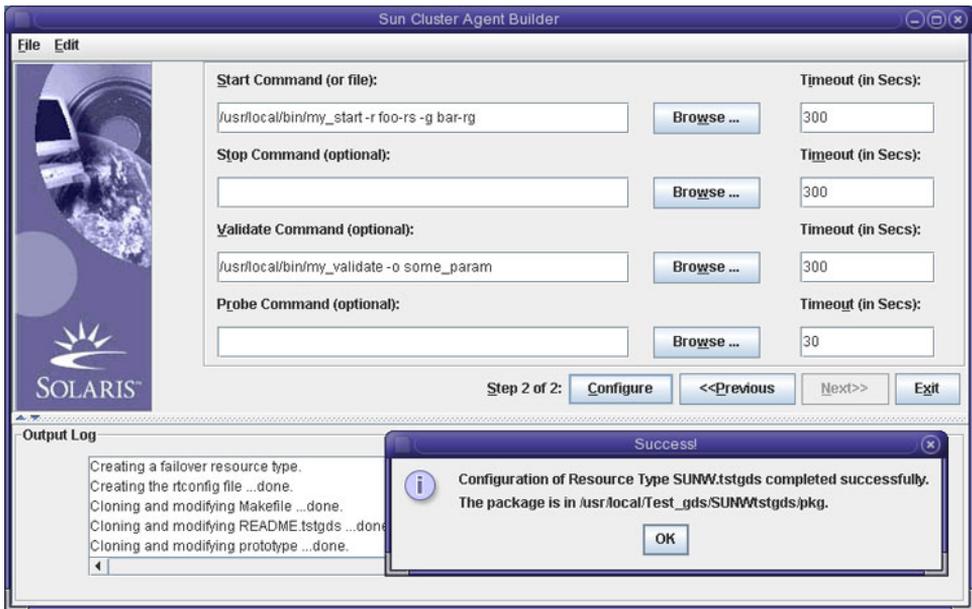


Figure 4.8 Completing the resource type definition using scdsbuilder

Integrating Your Application-Specific Logic

You use the fields in this second screen to provide the location of the programs (which can be compiled executables or scripts) and their associated arguments that will be used to start, stop, probe, and validate your data service when it is installed on the target cluster nodes. For each program, you can set a time limit on how long it can take for the program to complete. If the program does not complete within the allocated time period, then the resource is placed into a failed state, such as STOP_FAILED.

You are required to provide a value only for the start program. All the other programs are optional. Any programs specified must exit with a return code of zero

only when they have successfully completed their work. If they fail to perform their allotted task, they must return a value greater than 100. Values below that are used by the Solaris Cluster commands and have specific meanings (see the `intro(1CL)` man page).

The programs you assign to the start and stop commands must return successfully only when your target application has actually completed the relevant operation. If the stop command leaves the application under its control running, or not completely stopped, but the stop command returns successfully, then the cluster framework erroneously determines that it is safe to start the resource group on another cluster node. In some instances, particularly when the application uses a global file system, this outcome could result in data corruption because the two instances of the application could write to their data files in an uncontrolled fashion.

If no stop command is provided, the process tree that results from the start command is terminated using the `kill` command.

The `validate` command enables you to check that your application is correctly configured on all the potential nodes on which it can run. Again, if the program determines that your application is misconfigured, the `validate` program must exit with a nonzero exit code.

The capability to incorporate a probe command is one of the key benefits of using the Solaris Cluster framework. A probe command enables you to write a program that determines the health of your application. As an example, if you are writing a probe for a database, you could test whether it can execute basic SQL statements, such as creating or deleting a table, or adding or deleting a record. If you do not provide a probe script, then default methods are used instead.

For non-network-aware applications, the process-monitoring command `pmfadm` (see the `pmfadm(1M)` man page) monitors the process tree spawned by your start command. Only if all the processes have failed will the cluster framework attempt to restart the service. Therefore, if your service consists of multiple processes and only one process fails, then `pmfadm` will not recognize this fault unless it causes all the other processes to fail as well. Consequently, if you need to monitor your application with a higher degree of granularity, you must provide a custom fault probe.

If the application is network-aware, then the default probe tries to open the port listed in the `port_list` property. Because this is a simple probe, it makes no attempt to retrieve any data. Even if the default probe successfully opens the ports, that does not necessarily indicate overall application health.

In the preceding example, you would install the package generated by `scdsbuilder` on all your cluster nodes. You would then register the new resource type so that you could create new resources of this type. When the RGM is requested to create a resource, it calls the `validate` command: `/usr/local/bin/my_validate -o some_param`. If that command succeeds and you enable the resource, the RGM calls the `/usr/local/bin/my_start -r foo-rs -g bar-rg` command. In both

cases, the initial arguments are fixed, but you can modify them subsequently using the `clresource` command.

Resource Type Registration File

If you decide to write an agent from scratch using either the RMAPI or DSDL APIs, you must first describe the properties of your proposed resource type in a file known as the resource type registration (RTR) file. This file provides the RGM with details on which programs to call and which variables are required to control the particular application.

Example 4.18 shows an extract from the `SUNW.LogicalHostname` RTR file. As the example shows, all the programs for this resource type are located in the directory defined by `RT_BASEDIR`. The RTR file also defines programs that will, among other tasks, start, stop, and probe (`Monitor_start`) the logical IP address that the resource plumbs. These addresses are, in turn, defined in the `HostnameList` property.

The extension properties you define are all application-specific. They could, for example, refer to the location of the software binaries, that is, the application home directory. If a property has a default value, then you can define it in the RTR file to save your system administrator from having to override it each time he or she creates a resource of this type. Furthermore, you can place limits on what values certain properties can take and when they can be changed.

Example 4.18 Extract from the `SUNW.LogicalHostname` RTR File

The following text shows some of the key parts of the RTR file for the `SUNW.Logical-Hostname` resource type:

```
#
# Copyright 1998-2008 Sun Microsystems, Inc. All rights reserved.
# Use is subject to license terms.
#

#ident "@(#)SUNW.LogicalHostname 1.20 08/05/20 SMI"

# Registration information and Paramtable for HA Failover IPaddress
#
# NOTE: Keywords are case insensitive, i.e. users may use any
# capitalization style they wish
#

RESOURCE_TYPE ="LogicalHostname";
VENDOR_ID = SUNW;
RT_DESCRIPTION = "Logical Hostname Resource Type";
```

```

SYSDEFINED_TYPE = LOGICAL_HOSTNAME;

RT_VERSION = "3";
API_VERSION = 2;

INIT_NODES = RG_PRIMARYS;

RT_BASEDIR=/usr/cluster/lib/rgm/rt/hafoip;

FAILOVER = TRUE;

# To enable Global_zone_override
GLOBAL_ZONE = TRUE;

START                =          hafoip_start;
STOP                 =          hafoip_stop;

PRENET_START         =          hafoip_prenet_start;

VALIDATE             =          hafoip_validate;
UPDATE               =          hafoip_update;

MONITOR_START        =          hafoip_monitor_start;
MONITOR_STOP         =          hafoip_monitor_stop;
MONITOR_CHECK        =          hafoip_monitor_check;

PKGLIST = SUNWscu;

#
# Upgrade directives
#
#$upgrade
#$upgrade_from "1.0" anytime
#$upgrade_from "2" anytime

# The paramtable is a list of bracketed resource property declarations
# that come after the resource-type declarations
# The property-name declaration must be the first attribute
# after the open curly of a paramtable entry
#
# The Paramtable cannot contain TIMEOUT properties for methods
# that aren't in the RT
{
    PROPERTY = Start_timeout;
    MIN=360;
    DEFAULT=500;
}
.
.
.
# HostnameList: List of hostnames managed by this resource. All must be
# on the same subnet. If need > 1 subnet with a RG, create as many
# resources as there are subnets.
{
    PROPERTY = HostnameList;
    EXTENSION;
    STRINGARRAY;
    TUNABLE = AT_CREATION;
    DESCRIPTION = "List of hostnames this resource manages";
}
.
.
.

```

Resource Management API

The Resource Management API (RMAPI) is a set of low-level functions contained in the `libscha.so` library with both C and shell interfaces. All the function names provided by this interface are prefixed with `scha_`. The shell interfaces are listed in section 1HA of the Solaris Cluster manual pages.

The `ksh` scripts generated by the Agent Builder are built using these commands, so you can insert additional lines in this code where the comments indicate. However, for greater control over the logic imposed on your application you must write your application agent from scratch.

Data Service Development Library

The Data Service Development Library (DSDL) is a set of higher-level functions encapsulated in the `libdsdev.so` library that builds on the RMAPI functionality. This library can only be accessed using a C programming language interface. Consequently, it is potentially more time-consuming to write a complete application agent using this approach, although it does offer the greatest level of performance and flexibility.

If you used Agent Builder to create a resource type, you can customize it by inserting extra DSDL code where the comments indicate. Otherwise, you must write your agent from scratch.

All the function names provided by the library are prefixed with `scds_` and are documented in section 3HA of the Solaris Cluster manual pages. The NFS agent source code [NFSAgent] serves as a good example of how these APIs are used. Using the `nfs_svc_start.c` source as a specific example, the library is initialized with `scds_initialize()`. Resource and resource group names are then retrieved using `scds_get_resource_name()` and `scds_get_resource_group_name()` calls, respectively. Finally, the status of the resource is set by the RMAPI `scha_resource_setstatus()` call. Most of the coding effort involved with using these interfaces is consumed by the logic that describes how the agent should behave in various failure scenarios. For example, how many times should the agent attempt to restart the service before giving up and potentially failing over? What should the agent do in response to a network failure?

One advantage of using the GDS is that all the best practices for service behavior are already in the logic of the code that makes up the agent, saving you from re-creating that code.

Useful Utilities for Building Custom Data Services

The Solaris Cluster software comes with two programs that you will find very useful if you create your resource type from scratch: `hatimerun` (see the `hatimerun(1M)` man page) and `pmfadm`.

hatimerun Command

Throughout the `Start`, `Stop`, `Monitor_start`, and `Validate` methods of your resource type, you will need to run various programs to perform the required logic steps. Because your goal is high availability, you cannot wait for a program that might never respond or return, whether that program has gone into a loop or is unable to retrieve some important data from the network, disk, or other program. Consequently, you must place time constraints on the duration of the program's execution. This is the function of the `hatimerun` command. It enables you to execute a program under its control and set a limit on the time it can take to respond. If the program in question fails to respond in a timely fashion, it is terminated by default.

The `hatimerun` command also enables you to leave the program running asynchronously in the background, change the exit code returned after a timeout, or use a particular signal to terminate your program.

The most common usage of this command is in your probe commands or in the steps leading up to stopping or starting your application.

pmfadm Command

If you write a custom probe for your service, you decide what constitutes a healthy service. The criteria might include application-specific checks to determine if the data it is delivering to potential clients is valid or timely. If the application consists of multiple processes, you might want to check that each process is running, using the `ps` command. All of these tests combine to give you the best assessment of your application's current health. However, your probe is scheduled to make its checks only at regular intervals. Even though you can tune these checks to occur at shorter intervals, doing so results in a greater load on your system. Consequently, you must wait, on average, half the probe period before your probe detects a situation where your application has completely failed, meaning that all the processes have exited. Once again, this does not help much toward your goal of high availability.

The solution is to use `pmfadm`, the process-monitoring facility command. When you start your application under `pmfadm`, it monitors all the processes your application spawns to a level that you determine. By default, it monitors all the application's child processes. If they all exit, `pmfadm` immediately restarts your

application for you on the condition that it has not already exceeded a preset number of restarts within a certain time interval.

The most common usage of this command is in your start command to ensure that your key application processes are monitored and that complete failures are reacted to immediately.

libscho.st .so Library

Some applications store or make use of configuration information about the physical hostname of the server on which the application is running. Such applications will most likely fail when the application is placed in a resource group and moved between the nodes of a cluster. This failure occurs because calls to `uname` or `gethostbyname` produce different responses on the global zone of each cluster node. Oracle Application Server and the Oracle E-Business Suite are two examples of programs that risk such failures [LibHost].

To overcome this limitation, you use the `LD_PRELOAD` feature to enable the runtime linker to interpose the `libscho.st .so.1` library in the dynamic linking process. The following example shows how this is done. You can use the same construct within your resource `Start` or `Monitor_start` (`probe`) methods, as required.

Example 4.19 How to Use the `sclibhost .so.1` Library to Change the String Returned as the Hostname

Use the `uname` command to display the current hostname.

```
# uname -n
phys-winter1
```

Set the `LD_PRELOAD_32`, `LD_PRELOAD_64` and `SC_LHOSTNAME` environment variables, and then rerun the `uname` command.

```
# LD_PRELOAD_32=$LD_PRELOAD_32:/usr/cluster/lib/libscho.st .so.1
# LD_PRELOAD_64=$LD_PRELOAD_64:/usr/cluster/lib/64/libscho.st .so.1
# SC_LHOSTNAME=myhost
# export SC_LHOSTNAME LD_PRELOAD_32 LD_PRELOAD_64
# uname -n
myhost
```

Tuning and Troubleshooting

If you are running your applications on a Solaris Cluster system, then service availability is one of your main concerns. The two major causes of service migra-

tion are cluster-node failures and the resource probes detecting that an application is unhealthy.

The Solaris Cluster framework is responsible for detecting and reacting quickly to node failures. When all the cluster private interconnects fail, the default heartbeat settings allow only 10 seconds to elapse before a cluster reconfiguration is triggered and any affected service is restarted on one of the remaining cluster nodes. The application fault probe is responsible for detecting and reacting to an application that is either unhealthy or has failed. Each fault probe can be tuned to take the appropriate action based on the number of failures within a given time interval. Similarly, the resource type `Start` and `Stop` methods for the application are responsible for ensuring that the application is successfully started and stopped. The `Start`, `Stop`, and `Monitor_start` (fault probe) methods have a tunable timeout interval, which is the maximum duration for each operation. The fault probe also has a probe interval, which determines how often the probe is run.

When tuning the probe interval and any associated timeout, you must trade off the impact of the probe on system load and application performance. If you set the values too low, you risk false failover. If you set them too high, you risk waiting longer than necessary to detect a critical problem. Because all applications and workloads are different, the only realistic way to achieve an optimal setting is by thorough testing. You start with high values and gradually reduce them until lower values result in misdiagnosed problems. You then increase the values again to give the fault probe some scope of variability. Of course, all your testing must be performed under realistic workload conditions.

For the start and stop timeouts, your task is slightly easier because the system logs (in `/var/adm/messages`) state how much of the available timeout was used by each method to complete the relevant action. Thus, you can tune your start and stop timeouts such that the application takes no more than, for example, 50 percent of this value under normal conditions.

If it is not obvious why a resource is not functioning correctly, then you can often obtain additional debugging information to help diagnose the problem. The `/var/adm/messages` file is your main source of help. All your resources will log messages here. You can increase the amount of debugging information logged by ensuring that your `syslog.conf` file directs `daemon.debug` output to a suitable file.

Some resource types use flags for specifying the amount of debugging output they produce. Others written in DSDL code (see the section “Data Service Development Library”) rely on you specifying a log level from 1 to 9 (with 9 as the maximum) in a file named `/var/cluster/rgm/rt/resource-type-name/loglevel`.

This page intentionally left blank



Index

- A**
- Access log files, 384–385, 402
- Active membership, 44
- Active primary I/O paths, 60
- Adapter trap port, SNMP, 182
- add-trust subcommand
 - automating certificate exchange, 256
 - establishing trust, 254
- add_app_rg_args protection group property, 309
- add_app_rg_script protection group property, 309
- Address Resolution Protocol (ARP), 23
- Address short-listing, 95–96
- Administration server, Sun Java System, 383–384
- Administration, single point of, 153
- Administrative workload reduction, 154
- Administrator-initiated takeover, 227
- Advanced Agent Toolkit
 - extending GDS with, 202
 - supporting new applications, 207
- Advantages
 - of application-based replication, 232–233
 - of global networking service, 98
 - of Solaris Cluster, 1–5
- Affinities, resource group, 119–123
- Affinity_timeout property, 91
- Agent Builder command, 212–214
- Aggregated links, 25–26
- Alarm, disk utilization generating, 186
- Amnesia condition, 45, 51–52
- ANY_NODE, 114
- apacheset diskset
 - AVS initialization files for, 266–267
 - configuration of, 268–272
- app-pg protection group, 281
- Application
 - access to file systems, 68
 - executables, 75–76, 214
 - failure, GDS handling of, 205
 - resource groups. *See* Resource groups
- Application agents. *See also* Data services/
application agents
 - suitability criteria, 201–202
 - supporting new, 207
 - troubleshooting unhealthy, 220–221
- Application-based replication
 - advantages of, 232–233
 - disaster recovery and, 224
- Application data
 - automated migration and, 226

- Application data (*continued*)
 - backing up, 198, 199–200
 - consistency, 20, 273
 - HA containers and, 243
 - protecting, 260, 262–263
 - Application fault isolation, 152
 - Application-specific logic, 214–216
 - Architecture, Geographic Edition. *See also* Solaris Cluster Geographic Edition, features
 - automated vs. automatic migration, 226–227
 - N+1/paired topologies, 235–237
 - three-data-center. *See* Three-data-center (3DC) configuration
 - Architecture, PxFs, 69–71
 - Architecture, Solaris Cluster. *See also* Solaris Cluster, features
 - kernel/non-kernel components, 35–38
 - overview, 6
 - server configurations, 6–7
 - server-to-storage connectivity topologies, 8–12
 - storage connectivity, 7–8
 - Architecture, SPARC
 - Oracle VM server for. *See* Oracle VM server for SPARC
 - for quorum server, 50
 - for server configurations, 6
 - Architecture, zone-cluster
 - file systems, 157–158
 - membership monitoring, 155–156
 - overview, 154
 - security, 156–157
 - virtual node/membership, 155
 - ARP (Address Resolution Protocol), 23
 - ASM. *See* Automatic Storage Management (ASM)
 - ASYNc fence level, 282
 - Asynchronous replication
 - data consistency and, 19, 226
 - with Hitachi replication software, 282
 - manual intervention with, 18
 - MySQL capability for, 300
 - with storage-based replication, 21, 228
 - Asynchronous write operations, 72–73
 - At most once semantics, 40
 - automated, `RG_slm_type` property set to, 188–189
 - Automated vs. automatic service migration, 226–227
 - Automatic Storage Management (ASM), 66, 82, 106, 200
 - Availability, of Solaris Cluster software, 1–5
 - AVS. *See* StorageTek Availability Suite (AVS)
 - AVS snapshot files, 235, 266–267, 325
- B**
- Back-to-back private interconnects, 27–28
 - Backup
 - of application data, 199–200
 - clusters, 237–238, 315, 317
 - databases, 300
 - disaster recovery and, 224
 - HA containers, 243
 - heartbeat mechanisms, 259
 - hosts, 133
 - nodes, 10–11, 121
 - root disk, 198–199
 - servers, highly available, 200–201
 - The Beginners Guide to LDOMs: Understanding and Deploying Logical Domains*, 137
 - Benefits
 - of `clprivnet0` to Oracle RAC, 102–103
 - of shared QFS file system, 82
 - of third-party disaster recovery framework, 225
 - BIA (Business Impact Analysis), 225
 - Binaries
 - application, 74–76
 - GDS subclass, 208
 - installing Oracle, 372
 - storing Oracle, 82, 387
 - whole-root zone, 145
 - bind system
 - features of, 94–95
 - sticky/weighted policies and, 95–97
 - "Black-box" control
 - control domain and, 138
 - failover guest domains as, 245
 - HA containers as, 147–148, 243
 - of virtualized entities, 133–134
 - Boot protocol, 56
 - BrandZ framework, 146–147, 155
 - Buffered I/O, 79
 - Business Impact Analysis (BIA), 225

- C**
- C++ procedure calls, 38, 69
 - CAC. *See* Common Agent Container (CAC)
 - modules
 - cacoadm command
 - displaying CAC module status, 251
 - listing debugging filters, 328
 - listing trusted certificates, 253
 - SNMP MIBs and, 182
 - Cache coherence, for PxFs, 72
 - Cache Fusion, 102–103, 124
 - Cached I/O, 79
 - Callback methods, resource type, 109
 - Campus clusters
 - handling quorum in, 53–54
 - PxFs and, 77
 - stretching cluster node separation and, 28–29
 - Campus/metro cluster configuration, 53–54
 - CCR. *See* Cluster Configuration Repository (CCR)
 - cl_ccrad daemon, 126
 - cldevice command, 57–58
 - cldevicegroup command, 60, 62–63
 - Cleanup thread, for TCP connections, 93
 - cl_eventlogd daemon, 127
 - clexecd daemon, 126
 - CLI. *See* Command-line interface (CLI)
 - Client connection recovery. *See also* Recovery, 98–99
 - Client push communication technique, 73
 - cl_pnmd daemon
 - function of, 127
 - public networking and, 23
 - clprivnet software, 37
 - clprivnet0 virtual network interface
 - for aggregation of cluster interfaces, 101–102
 - benefits to Oracle RAC, 102–103
 - clquorumserver command, 50
 - clresourcegroup command
 - bringing resource groups online, 115
 - changing RGM resource group information, 43
 - for resource group property settings, 188–189
 - showing resource group status, 250–251
 - Cluster backup. *See also* Backup
 - application data backup, 199–200
 - HA backup servers, 200–201
 - overview, 198
 - root disk backup, 198–199
 - Cluster brand zones, 155–156
 - cluster command, 58
 - Cluster Configuration Repository (CCR)
 - configuration information in, 324
 - configuration information storage, 252–253
 - dual-partition upgrades and, 193
 - file consistency of, 43
 - functions of, 42–43
 - manual editing of files, 44
 - zone-cluster file systems and, 158
 - Cluster device fencing options, 58
 - Cluster events, SNMP MIBs for, 182
 - Cluster file systems
 - application binaries/data/logs, 75–76
 - campus cluster and PxFs, 77
 - Mount subsystem, 74
 - overview/application access, 68
 - protecting Oracle RAC databases and, 313
 - PxFs architecture, 69–70
 - PxFs at file system level, 70–71
 - PxFs cache coherence, 72
 - PxFs compared to NFS, 77–79
 - PxFs files/HA service, 71–72
 - PxFs I/O support, 72–73
 - PxFs performance, 76–77
 - zone clusters supporting, 158
 - Cluster heartbeat messages. *See* Heartbeat messages
 - Cluster kernel/non-kernel components, 35–38
 - Cluster management
 - command-line interface for, 173–174
 - with Geographic Edition software, 324–325
 - overview, 173
 - role-based access control, 178–180
 - Solaris Cluster Manager GUI for, 175–177
 - Solaris Cluster Manager wizards, 177–178
 - Cluster membership
 - determining, 44–45
 - majority voting. *See* Majority voting
 - split-brain condition and, 44
 - Cluster Membership Monitor (CMM)
 - functions of, 44–45
 - implementation of, 45

- Cluster Membership Monitor (*continued*)
 - maintaining cluster integrity, 100
 - programming interface, 54–55
 - rgmd daemon and, 107
- Cluster monitoring
 - overview, 180
 - SNMP management information bases, 182–183
 - Sun Management Center integration, 181
- Cluster node membership, 55
- Cluster node separation, 28–29
- Cluster node votes, 45–46
- Cluster nodes
 - backup, 10–11, 121
 - configuring LUNs to, 283
 - configuring Solaris Cluster software on
 - first, 344–352
 - second, 352–356
 - configuring/validating, 172
 - defined, 132–133
 - I/O domains used as, 138–139
 - installing Oracle 11g software on both, 358–362
 - installing Oracle Solaris OS on, 163–166
 - time synchronization across, 172–173
 - troubleshooting failures, 220–221
 - upgrading with zones, 196–198
 - virtual, 155–156
- Cluster objects, 183
- Cluster private interconnect
 - benefits to Oracle RAC, 102–103
 - clprivnet0 virtual network interface, 101–102
 - configuration guidelines, 105
 - InfiniBand and, 104
 - overview/heartbeats, 99
 - resilience/protocols/TCP/IP, 103–104
 - subnet allocation/network ports, 104–105
 - TCP/IP and, 103–104
 - topology, 99–100
 - traffic, 100–101
 - zone clusters supporting, 160
- Cluster protocol, 56
- Cluster public networking, 86–88
- Clustered pair storage topology, 9–10
- cluster_interconnects parameter, 102
- Cluster(s)
 - backup, 237–238, 315, 317
 - creating multiple distinct, 134–136
 - creating partnerships, 254–257
 - creating single-node system, 335–339
 - creating trust, 253–254, 323
 - defined, 133, 235
 - patching/upgrading. *See* Patching/upgrading
 - reconfiguration process, 55–57
 - resilience to network failure, 103
- CMM. *See* Cluster Membership Monitor (CMM)
- Code, binding to network address, 88
- Command device LUNs, 283
- Command-line interface (CLI)
 - for cluster management, 173–174
 - Geographic Edition, 325
 - implementation, 252
- Commands, Solaris Cluster CLI, 174
- Common Agent Container (CAC) modules
 - debugging, 327–329
 - for Geographic Edition software, 249–251
- Common Object Request Broker
 - Architecture (CORBA) model, 38–39
- com.sun.cluster.geocontrol CAC module, 251
- com.sun.cluster.notifier CAC module, 252
- Configuration
 - control, 42–44
 - guidelines, 105
 - information, storage of, 252–253
 - server, 6–7
 - type, 342
- configuration_file protection group property, 309
- Configuring trust and partnership, 253–257, 323
- Connection topologies, Geographic Edition, 235–237
- Consistency. *See* Data consistency
- Consistency group ID (ctgid) property
 - data consistency with, 20
 - matching resource/protection group, 287
 - with multiple device groups, 282
- Constructs, data service, 107

- Containers, Oracle Solaris. *See* Oracle Solaris Zones (Containers)
- `con_threshold` property, 93–94
- Control domains, 137–138
- CoolThreads, Oracle Solaris, 136, 140, 244
- CORBA (Common Object Request Broker Architecture) model, 38–39
- Core patch, Solaris Cluster software, 343
- Cost-effectiveness, of HA framework, 5
- `create_config_args` replication component property, 311
- `create_config_script` protection group property, 309, 312
- `ctgid` property. *See* Consistency group ID (`ctgid`) property
- Custom heartbeat mechanisms, 259
- Custom installation, of Sun Java System Web server software, 383
- `cznetd` daemon, 127
- D**
- Daemon processes. *See also specific daemons*
 - list of, 126–129
 - in Oracle Clusterware framework, 124
 - public network monitoring, 23
 - quorum server, 50–51
- Data, application. *See* Application data
- Data centers
 - performance impact of separating, 30–34
 - quorum devices in, 53–54
- Data consistency
 - after rolling failures, 19–20
 - of applications, 20, 273
 - disaster recovery and, 226
 - with Hitachi replication software, 282
 - in host-based replication between sites, 241
 - of ZFS file system, 234–235
- DATA fence level, 282
- Data Guard. *See* Oracle Data Guard software
- Data, in cluster file systems, 75–76
- Data Link Provider Interface (DLPI), 99, 101
- Data protection. *See also* Replication, for data protection
 - disksets/disk groups, 14–17
 - overview, 13–14
 - storage-based replication for, 18–21
 - zpools, 17
- Data redundancy, 13
- Data replication. *See also* Replication, for data protection
 - creating additional technology for, 329–330
 - effect of starting/stopping protection groups, 314–315
 - modules. *See* Replication modules
 - protection groups without, 261–262
- Data replication resource group, 308
- Data Service Development Library (DSDL)
 - application integration and, 202
 - creating new resource types and, 218
 - `scdsbuilder` command and, 212
- Data services/application agents
 - constructs/`rgmd` daemon, 107
 - daemon processes, 126–129
 - defined, 37
 - overview, 105–106
 - parallel services, 123–126
 - resource groups, 117–123
 - resource types, 108–111
 - resources, 112–116
- Data volume LUNs, 283
- Database file location, for Oracle 11g, 366
- `dbca` configuration process, 364–366
- DCS. *See* Device Configuration Service (DCS)
- Debugging information, 327–329
- Dedicated cluster model, 153
- Dense wave division multiplexors (DWDMs), 28–29, 31
- Dependencies, resource, 113–115
 - `/dev/did` device vs. `/dev/global` device, 62–64
 - `/dev/did` DID devices, 66
 - `/dev/global` Global devices, 67
 - `/dev/md/disket` Global devices, 67
 - `/dev/vx/rdisk/disk_group` Global devices, 67
- `dev_group` property, 290
- Device Configuration Service (DCS), 57, 70
- Device group resources, scalable, 125, 126
- Device groups
 - creating SRDF, 274–277
 - data service and, 117
 - `rawdsk`-type, 64
 - Solaris Cluster to match SRDF, 277
 - working with, 60–63
- Device ID (DID) implementation, 61–64

- Device namespace, 64
 - `device_group` property, 272
 - Devices
 - choosing correct, 66–67
 - device ID, 61–64
 - disk path monitoring, 65–66
 - global, 59–60
 - namespace for, 64
 - overview, 59
 - `DG_or_CG` property, 280
 - DHCP. *See* Dynamic Host Configuration Protocol (DHCP)
 - Direct access I/O, 81
 - Direct I/O, 80
 - Directed acrylic graph, 115
 - Disabled resources, 112
 - Disaster recovery solution
 - in 3DC configuration, 238
 - choosing appropriate, 224–225
 - connection topologies and, 236–237
 - features of, 4
 - multi-node, multi-array cluster in, 242
 - reasons for, 223–224
 - third-party framework, 225
 - Disk failfast driver, 58
 - Disk fencing, 57–58, 159
 - Disk groups
 - correct configuration of, 125
 - data protection and, 14–15
 - Disk Path Monitoring (DPM) facility, 65–66
 - Disks, 158–159
 - Disksets (metasets)
 - adding third-site mediator to, 16–17
 - correct configuration of, 125
 - multi-owner, 15
 - single-owner, 14
 - state replica majority on, 15–16
 - Distributed lock manager (DLM), 123–124
 - `dladm` command, 25–26
 - DLPI (Data Link Provider Interface), 99, 101
 - Domains
 - creating multiple distinct clusters, 134–136
 - I/O, 138–139
 - roles for, 137–138
 - `domino` mode, 273
 - DPM (Disk Path Monitoring) facility, 65–66
 - DSDL. *See* Data Service Development Library (DSDL)
 - Dual-partition upgrades, 192–194
 - DWDMs (dense wave division multiplexors), 28–29, 31
 - Dynamic Host Configuration Protocol (DHCP)
 - exclusive/shared networking options and, 145
 - non-running in Oracle Solaris zones, 146
 - Dynamic system domains. *See also* Domains features of, 134–136
 - Geographic Edition with, 242
- E**
- EMC PowerPath, 13
 - EMC Symmetrix Remote Data Facility (SRDF)
 - configuration, 274–277
 - data protection and, 13, 18, 228–229
 - Geographic Edition resource group configuration, 277–280
 - monitoring, 280–281
 - overview, 273–274
 - Enabled resources, 112
 - Environments, virtualized, 232–233
 - Establishing trust, between clusters, 253–254, 323
 - `/etc/horcm.conf` file, 283–284
 - Ethernet adapter cards, 24–25
 - Event propagation, 252
 - Exactly once semantics, 40
 - Example implementations, of Solaris Cluster services. *See* Solaris Cluster implementations
 - Exclusive IP networking options, 145
 - Explicit binding, 94
 - External firewalls, 169
 - `External_Dependency_Allowed` protection group property
 - null data replication type and, 312
 - protecting Oracle RAC databases and, 313
- F**
- Failfast driver, disk, 58
 - Failfast panic, 57
 - `failfastd` daemon, 127

- Failover applications
 - cluster file systems and, 75, 76
 - cluster-wide zones and, 152
 - defined, 37–38
 - dual-partition upgrades and, 194
 - failover file system supporting, 68
 - private interconnect and, 101
 - Failover data service, 105
 - Failover file systems
 - disaster recovery protection and, 243
 - function of, 68
 - highly available local, 84–85
 - for Oracle database creation, 362
 - setting up/protecting, 287–289
 - Failover guest domains
 - constructing, 141–142
 - data services and, 106
 - using Geographic Edition with, 245
 - Failover resource groups, 97, 98, 110, 117, 119, 120
 - Failover resources, 116
 - Failover times, 45, 134
 - `Failover_mode` property, 205–206
 - Failure notification, 259
 - `failure-policy`, 141–142
 - Failure(s)
 - GIF node, 98–99
 - node, tuning/troubleshooting, 220–221
 - operating system, 2–3, 4
 - private network, 103
 - single point of, 6–7
 - Fair share scheduler (FSS), 184–185
 - FAN. *See* Oracle Fast Application Notification (FAN)
 - Fast file system recovery, 80
 - Fault probes
 - disaster recovery and, 227
 - function of, 5
 - for unhealthy applications, 220–221
 - Fault-tolerant systems vs. Solaris Cluster solution, 3
 - faulted paths, 100–101
 - Features
 - of Geographic Edition. *See* Solaris Cluster Geographic Edition, features
 - of Solaris Cluster. *See* Solaris Cluster, features
 - `fence_level` property, 282
 - Fencing, of disks/RAID devices, 159
 - Fencing subsystem, storage device
 - disk failfast driver, 58
 - disk fencing, 57–58
 - NAS device fencing, 59
 - overview, 57
 - Fibre Channel storage devices
 - creating metro clusters and, 28–29
 - in Solaris Cluster configurations, 7
 - File consistency, 43–44
 - File support, FxFS architecture for, 71
 - File systems
 - cluster. *See* Cluster file systems
 - failover. *See* Failover file systems
 - highly available, 157–158
 - local, 75–76, 84–85, 157
 - network. *See* Network file systems (NFSs)
 - Oracle Solaris ZFS. *See* Oracle Solaris ZFS file system
 - overview, 68
 - protecting with replication, 233–235
 - scalable mount-point resources for, 125
 - Sun QFS, 79–84
 - zone-cluster, 157–158
 - Fine-grained control, 133–134
 - Firewall(s)
 - external, 169
 - rules, 87–88
 - Flag files, 329
 - Flexibility, with application-based replication, 232–233
 - `forcedirectio` mount option, 74
 - Framework, high-availability. *See* High-availability (HA) framework
 - `FROM_RG_AFFINITIES`, 114–115
 - FSS (fair share scheduler), 184–185
 - `fssnap` option
 - for application data backup, 200
 - for root disk backup, 198–199
- ## G
- `gc_hb_heartbeat`-name format, 253
 - Generic Data Service (GDS)
 - creating new resource types, 202, 203–207
 - creating subclass of, 208–212
 - `scdsbuilder` command and, 212–214

- generic_affinity property
 - global networking and, 91
 - outgoing_connection property and, 97
 - geo-clustername resource, 249–251
 - geo-clusterstate resource group, 250–251
 - geo-failovercontrol resource, 249–251
 - geo-infrastructure resource group, 249–251
 - geoadm command, 325
 - Geographic Edition. *See* Solaris Cluster Geographic Edition
 - Geographic Edition clusters
 - creating trust between, 253–254
 - partnerships between, 254–257
 - geohb command, 325
 - geopg command
 - features of, 325
 - protection groups and, 260–263
 - for switchovers, 315–316
 - for takeovers, 319
 - geops command
 - displaying heartbeat information, 257–258
 - establishing trust between installations, 254–257
 - features of, 325
 - GIFs. *See* Global interfaces (GIFs)
 - GIN. *See* Global interface node (GIF node/GIN)
 - Global cluster
 - CCR storing, 42–44
 - features of, 133
 - Global-cluster membership, 55–56
 - Global-cluster non-voting nodes
 - creating scalable web services in. *See* Sun Java System Web server service, creating
 - defined, 133
 - Oracle Solaris zones using, 143–144, 148–149
 - using Geographic Edition with, 244
 - Global-cluster voting nodes, 133
 - Global devices
 - overview, 59–60
 - primary/secondary I/O paths, 60
 - PxFS file system mounted on, 70
 - /global/.devices/node@X directory, 64
 - Global fencing policy, 58
 - Global interface node (GIF node/GIN)
 - failure, 98–99
 - global networking and, 89–91
 - Global interfaces (GIFs)
 - global networking and. *See also* Global networking service
 - global networking service and, 89
 - for scalable services, 86
 - Global IP addresses. *See* Global interfaces (GIFs)
 - Global networking service
 - advantages of, 98
 - client connection recovery, 98–99
 - deterministic distribution, 92–94
 - generic affinity/affinity timeout, 91
 - IPsec security association failover, 97–98
 - outgoing connection support, 97
 - overview, 89
 - packet distribution mechanisms, 89–91
 - round_robin property and, 91–92
 - scalable service support for SCTP, 94–97
 - Global Networking subsystem, 37
 - Global zones
 - creating zone clusters from, 388–390
 - for quorum server, 51
 - GNS (Group Name Services), 274
 - Graphical user interface (GUI)
 - Geographic Edition Manager, 325–326
 - Sun Management Center, 181
 - Graphical user interface (GUI), Solaris Cluster Manager
 - for cluster management, 175–177
 - implementation, 252
 - installing Solaris Cluster software and, 170
 - scdsbuilder, 202, 212–214
 - Grids, 3
 - Group Name Services (GNS), 274
 - Guest domains
 - configuring, 139–140
 - failover, 141–142
 - role of, 137
 - using Geographic Edition with, 245
 - virtualized I/O and, 140–141
 - Guidelines, configuration, 105
- ## H
- HA. *See* High-availability (HA)
 - HA-Oracle 11g Release 1 database. *See* Oracle 11g Release 1 database, creating

- ha_mysql_config file, 302–303
 - hatimerun command, 219
 - HDLM (Hitachi Dynamic Link Manager), 13
 - Heartbeat messages
 - DLPI and, 99
 - private networking and, 26–28
 - Heartbeat messages, Geographic Edition
 - failure notification, 259
 - heartbeat module, 258–259
 - monitoring, 326
 - overview, 257–258
 - replication components and, 264–265
 - Heartbeat module, Geographic Edition, 258–259
 - heartbeat_quantum property, 100–101
 - heartbeat_timeout property, 100–101
 - High-availability (HA) containers. *See also* Oracle Solaris Zones (Containers)
 - for Oracle Solaris zones, 147–148
 - upgrading, 196–198
 - using Geographic Edition with, 243
 - High-availability (HA) framework
 - disaster recovery vs., 4
 - fault probes for, 5
 - PxFS server subsystem as, 71–72
 - services of, 37, 39–42
 - Solaris Cluster solution vs., 3
 - High-availability (HA) invocations, 40–42
 - High capacity, of Sun QFS file systems, 80
 - High performance computing (HPC) vs. Solaris Cluster solution, 3
 - Highly available applications, 201–202
 - Highly available backup servers, 200–201
 - Highly available file systems, 35, 157–158
 - Highly available local (failover) file systems, 84–85
 - Hitachi Data Systems Universal Replicator
 - adding resource groups, 287–289
 - communications failure and, 235
 - configuring files, 283–284
 - creating protection groups, 284–286
 - data consistency with, 19–20
 - data protection and, 13–14, 18, 228–229
 - journal volumes and, 283
 - monitoring, 289–291
 - overview, 282–283
 - resource group configuration, 289
 - Hitachi Dynamic Link Manager (HDLM), 13
 - Hitachi TrueCopy
 - adding resource groups, 287–289
 - CAC debug output for, 328–329
 - configuring files, 283–284
 - creating protection groups, 284–286
 - data protection and, 13–14, 18
 - monitoring, 289–291
 - overview, 282–283
 - resource group configuration, 289
 - horcmd daemon process, 283–284
 - Host-based mirroring
 - with host-based replication, 231, 241–242
 - storage-based replication vs., 18–21
 - Host-based replication
 - disaster recovery and, 224, 228
 - file system protection with, 233–235
 - with host-based mirroring, 241–242
 - with storage-based replication, 239–241
 - StorageTek AVS, 230–231
 - Hostnames
 - changing string returned as, 220
 - Data Guard Broker CLI and, 293–294
 - matching cluster names, 251
 - for Sun Java System administration server, 384
 - Hosts
 - configuring IPMP group and, 22
 - data protection and, 13
 - logical, 86–88, 201
 - shared QFS file systems and, 81–83
 - Solaris, 132–133, 150–151
 - storage connectivity and, 7
 - HPC (high performance computing), 3
 - Hypervisor, 137
- ## I
- I/O connectivity, 65–66
 - I/O domains
 - as cluster nodes, 138–139
 - role of, 137
 - using Geographic Edition with, 245
 - I/O latency, 30–34
 - I/O paths, 60, 63
 - icmp-echo packets, 22
 - ifconfig commands, 161
 - ifconfig_proxy_serverd daemon, 127
 - IKE (Internet Key Exchange), 98

- Implementing Solaris Cluster services. *See* Solaris Cluster implementations
 - Implicit binding, 94
 - InfiniBand, 104
 - Infrastructure
 - site for secondary IT, 224
 - software, Geographic Edition. *See* Software infrastructure, Geographic Edition
 - Initialization files, AVS, 266–267
 - `in.mpathd` daemon, 22–23
 - Installation
 - of Geographic Edition software, 321–324
 - of Oracle 11g software, 358–362
 - of Oracle Solaris OS, 163
 - of Solaris Cluster software, 339–343
 - of StorageTek Availability Suite, 399–400
 - of Sun Java System Web server software, 382–384
 - Installation directory, for Sun Java System, 382
 - Installation menu, for Solaris 10 OS, 332–333
 - Installation screen
 - for Oracle 11g software, 360
 - for Solaris Cluster software, 171, 340
 - `installer` program
 - for Geographic Edition software, 322
 - installation process for, 169–172
 - installing Oracle Solaris OS, 163
 - updating alternate boot environment, 195
 - Instruction architecture, for server configurations, 6
 - Interconnect, cluster private. *See* Cluster private interconnect
 - Interface bonding, 25
 - Internet Key Exchange (IKE), 98
 - Internet Protocol Suite (TCP/IP), 103–104
 - Intracenter host-based replication, 239–241
 - Intracenter storage-based replication, 239–241
 - Inventory location, Oracle 11g software, 361
 - Invocations
 - to HA services, 40–42
 - ORB and, 38
 - `ioctl` system calls, 159
 - IP addresses
 - associated with hostnames, 251
 - global. *See* Global IP addresses
 - local. *See* Logical IP addresses
 - private interconnect and, 160
 - public networking and, 160–162
 - IP exclusive option, 162
 - IP network multipathing (IPMP)
 - configuring for public networks, 22–23
 - public networking and, 21
 - script for static probe targets, 22–23
 - server configurations and, 6–7
 - zone-cluster public networking and, 160
 - IP Security Architecture (IPsec)
 - ensuring replicated data security, 254
 - SA failover for scalable services, 97–98
 - iSCSI storage devices
 - software quorum and, 50
 - storage connectivity and, 7–8
 - Isolation
 - application fault, 152
 - security, 151–152
- J**
- Java SE Software Development Kit upgrade, 340
 - Journal devices, 19–20
 - Journal volumes (LUNs), 283
 - Jumbo frames, 24–25
 - JumpStart Enterprise Toolkit (JET), 169–170
 - "Just A Bunch Of Disks" (JOBOD) arrays, 14
- K**
- Kernel components, of Solaris software, 35–38
 - `ksh` registration script
 - Resource Management API and, 218
 - `scdsbuilder` command and, 212, 214
- L**
- Latency
 - I/O, 30–34
 - network, 104
 - `LB_STICKY/LB_STICKY_WILD` policy
 - binding with, 95–96
 - global networking and, 90–92
 - `outgoing_connection` property and, 97
 - `LB_WEIGHTED` policy
 - binding with, 95–96
 - global networking and, 89
 - `libdsdev.so` library, 218
 - `libscha.so` library, 218

- libscho`st.so.1` library, 220
 - Link aggregation, 24–26
 - Load-balancing policies
 - of `clprivnet0`, 102–103
 - outgoing_connection property and, 97
 - of resource groups, 120–121
 - SCTP feature, 94
 - Local addresses, 95–96
 - Local file systems
 - cluster file system vs., 75–76
 - highly available (failover), 84–85
 - zone clusters supporting, 157
 - local_mac_address?`=true`, 22
 - local_database_name property, 295
 - local_db_service_name property, 296
 - local_disk_queue property, 272
 - LOCAL_NODE resource dependency, 114
 - local_rac_proxy_svr_rg_name property, 295
 - local_service_password replication component property, 311
 - Location choice, for Oracle database files, 366
 - lockfs option, 198–199
 - Log files, 75–76
 - Logical Domains. *See also* Oracle VM server for SPARC, 244–245
 - Logical hostnames
 - initialization files and, 266
 - MySQL replication and, 300
 - resource group configuration and, 268–272
 - Logical IP addresses
 - cluster public networking and, 86–88
 - for Geographic Edition installation, 322
 - global networking service. *See* Global networking service
 - Logical LUNs, 229
 - Logical standby database, 291
 - Logical (virtual) hosts
 - configuring firewall, 87–88
 - creating resource group with, 86–87
 - HA backup servers and, 201
 - Lookups, in Name Server, 38
 - Loss of "state," 3
- M**
- MAC-to-IP address mapping, 23
 - Majority voting
 - amnesia condition, 51–52
 - campus/metro cluster configuration, 53–54
 - membership programming interface, 54–55
 - monitoring feature, 53
 - network-attached storage, 51
 - overview, 45–47
 - quorum devices, 47–48, 54
 - quorum disk devices, 48–49
 - quorum server, 50–51
 - software quorum, 49–50
 - uneven cluster partitions, 52–53
 - Managed state, of resource groups, 117
 - Management information bases (MIBs)
 - for cluster events, 182
 - for cluster objects, 183
 - Mandatory locks, 82
 - Manual alteration, of shared storage reservations, 58
 - Manual editing, of CCR files, 44
 - Manual update, of shared QFS file system, 234
 - Master (primary) database, 300
 - Master servers, highly available, 200–201
 - Maximum availability mode, 292
 - Maximum performance mode, 292–293
 - Maximum protection mode, 292
 - Maximum Tolerable Period of Disruption (MTPD), 225
 - mdmonitord daemon, 127
 - Mediators
 - adding to disksets, 16–17
 - state replica majority and, 15–16
 - Membership
 - cluster, 44–45
 - monitoring, 155–156
 - zone-cluster, 155
 - Membership programming interface, 54–55
 - Membership subsystem, 37
 - Metadata storage, 80–81
 - Metadevices (RAID devices)
 - data protection and, 14
 - zone clusters supporting access to, 158–159
 - Metasets. *See* Disksets (metasets)
 - Metropolitan (metro) clusters
 - creating, 28–29
 - handling quorum in, 53–54
 - MHI`CENFAILFFAST` parameter, 57

- MIBs. *See* Management information bases (MIBs)
 - Migration, of services
 - automated vs. automatic, 226–227
 - switchover as, 315–316
 - takeover as, 317–319
 - Minimal performance overhead, 144–145
 - Minimization, of operating system, 166–167
 - Mirror consistency, 282
 - mode_sc option, 139
 - Modules, replication
 - creating additional technology for, 329–330
 - debugging information for, 327–328
 - implementation, 252
 - Mount subsystem, for PxFs, 74
 - MPxIO, 7, 13
 - MTPD (Maximum Tolerable Period of Disruption), 225
 - Multi-master data service, 105
 - Multi-master resources, 116
 - Multi-owner (shared) disksets, 15
 - Multi-threaded/single-threaded write performance, 32–33
 - Multihoming, 94
 - Multipathing. *See also* IP network multipathing (IPMP)
 - data protection and, 13–14
 - with MPxIO, 7, 13
 - MySQL replication
 - configuration, 301–306
 - creating additional replication technology for, 329–330
 - overview, 300
 - mysql_config_file file, 301–302
 - mysql_geo_config file, 303–304
 - mysql_geo_register command, 305
- N**
- N+1 storage topology
 - features of, 10–11
 - Geographic Edition, 235–236
 - Name Server
 - function of, 36
 - HA service and, 40
 - ORB and, 38
 - Negative resource group affinity, 118
 - Network-attached storage (NAS) devices
 - fencing, 59
 - Pair+M topology and, 12
 - as quorum devices, 51
 - Network communications
 - securing, 169
 - zone clusters supporting, 160–162
 - Network file system (NFS)
 - NAS device fencing in, 59
 - non-running in Oracle Solaris zones, 146
 - public networking and, 22
 - PxFs compared to, 77–79
 - zone clusters supporting, 158
 - Network interface cards (NICs)
 - clprivnet0 support of, 102
 - creating two tagged VLANs on, 24
 - exclusive/shared networking options and, 145
 - public networking and, 21
 - server configurations and, 6–7
 - zone-cluster public networking and, 160–162
 - Network performance
 - impact of separating data centers, 30–34
 - improving with link aggregation/jumbo frames, 24–26
 - Network ports, 105
 - Network Time Protocol (NTP), 101, 172–173
 - NEVER fence level, 282
 - New resource types. *See* Resource types, creating
 - NFS. *See* Network file system (NFS)
 - NIC. *See* Network interface cards (NICs)
 - No space errors, 73
 - Node pair protocol, 56
 - nodelist property
 - function of, 118
 - resource groups in different zones and, 123
 - Nodes. *See* Cluster nodes
 - nofencing setting, 57
 - Non-global zones. *See also* Global-cluster non-voting nodes
 - configuring, 143–144
 - defined, 51
 - installing Sun Java System Web server software from, 382–384

Non-kernel components, of Solaris software, 35–38

Nonshared (single-owner) disksets, 14, 16–17

`Notification_EmailAddr`s property, 259

NTP. *See* Network Time Protocol (NTP)

Null (none) data replication type, 312

O

Object Request Broker (ORB) subsystem, 37, 38–39

Offline restart dependency, of resources, 112

Oracle 11g Release 1 database, creating
 completing database setup, 367–372
 configuration information for, 364–366
 installing software on both nodes, 358–362
 preparing for database creation, 362–364

Oracle Cluster Registry (OCR)
 in Clusterware framework, 124–125
 replication of storage, 313

Oracle Clusterware framework
 availability management with, 124–125
`clprivnet0` and, 102

Oracle Data Guard Broker
 protection group configurations, 295
 for replication, 293–294

Oracle Data Guard software
 configuration, 294–298
 Oracle RAC databases and, 293–294
 overview, 291–293
 resource group configuration/monitoring, 299

Oracle Enterprise Manager database control service
 choosing not to configure, 365
 creating, 367–372
`SUNW.gds` resource for, 358–360

Oracle Fast Application Notification (FAN), 126

Oracle RAC wizard, 178

Oracle Real Application Clusters (RAC)
 availability management integration
 and, 124–125
 benefits of `clprivnet0`, 102–103
 data protection and, 15
 databases, 293–294, 313
 disk fencing and, 58
 read-only restriction on, 20

`ORACLE_SID`, 149, 178, 295, 364, 368–371, 391, 393

Oracle Solaris Cluster. *See* Solaris Cluster

Oracle Solaris Cluster Geographic Edition.
See Solaris Cluster Geographic Edition

Oracle Solaris operating system (OS). *See* Solaris operating system (OS)

Oracle Solaris ZFS file system
 as failover file system, 84, 85, 362
 protecting with replication, 234–235
 snapshot feature, 199, 200, 224
 storage pool RAID protection, 17

Oracle Solaris Zones (Containers)
 branded zone framework in, 146–147
 global-cluster non-voting nodes, 148–149, 244

HA containers, 147–148, 243

IP exclusive/shared networking/root directory, 145

minimal performance overhead, 144–145
 overview, 143–144

patch process/non-running services, 146
 upgrading nodes with, 196–198
 using Geographic Edition with, 243–244

Oracle VM server for SPARC, 136–138

Oracle VM VirtualBox software
 creating single-node cluster system, 335–339

installing Solaris 10 OS, 331–335

ORB (Object Request Broker) subsystem, 37, 38–39

`outgoing_connection` property, 97

P

Packet distribution mechanisms, 89–91

Paged I/O, 79

Paired topology, 235–236

Pair+M storage topology, 11–12, 52, 100

Panic, failfast, 57

Parallel applications
 Clusterware and, 124–125
 overview, 123–124
 restriction on, 20

Solaris Cluster resources, 125–126

Parallel data service, 105

Partitions, dual-partition upgrades and, 192–194

- Partnerships
 - creating, 235
 - between Geographic Edition clusters, 254–257, 395, 398
 - replication components and, 264–265
- Passive secondary I/O paths, 60
- Patching/upgrading
 - of Geographic Edition software, 324
 - HA container availability and, 147–148
 - for Oracle Solaris zones, 146
 - overview, 189
 - upgrade methods, 190–195
 - upgrading nodes with zones, 196–198
- Per-node resource properties, 115–116
- Performance
 - network, 24–26, 30–34
 - overhead, 144–145
 - PxFS, 76–77
- Persistent Group Reservation emulation (PGR_e) extension, 48–49
- pg_protection-group-name format, 253
- Physical standby database
 - of Data Guard Broker configuration, 296
 - defined, 291
- "Ping-pong" prevention, 118
- Pingpong_interval property, 205–207
- Plex, host-based mirroring and, 18
- pmfadm command, 219–220
- pmmmd daemon, 127
- Ports, network, 105
- Positive resource group affinity, 118
- POSIX semantics, 72–73
- Prerequisite checks, for software
 - installation, 361
- Primary databases, 291
- Primary I/O paths
 - to device groups, 63
 - to global devices, 60
- Primary node failure, PxFS, 41–42
- Primary providers, for HA services, 39–40
- Primary sites
 - automated vs. automatic migration and, 226–227
 - combined host-based mirroring/replication and, 241–242
 - connection topologies and, 237
 - disaster recovery and, 4, 223–225
 - domino mode and, 273
 - failover guest domains and, 245
 - host-based replication and, 231
 - maximum performance mode and, 293
 - mirror consistency in, 282
 - mixing storage-based/host-based replication and, 238–240
 - Oracle Data Guard configuration and, 294–296
 - protection groups at, 261, 315
 - in shared QFS file system, 233–234
 - storage-based replication and, 229, 238
 - takeovers and, 317
- Private interconnect. *See* Cluster private interconnect
- Private interconnect heartbeats. *See* Heartbeat messages
- Private interconnect topology, 99–100
- Private interconnect traffic, 100–101
- Private network failure, 103
- Private networking
 - features of, 26–28
 - protocols, 103
 - zone clusters supporting, 160
- Probe targets, 22–23
- Probe_command property
 - function with new program, 214–216
 - of GDS agent, 203
 - of GDS subclass, 208–211
- Process membership, 55
- Process up/zone down, 55
- Processor pool, 184–185
- Programming interface, CMM, 54–55
- Protection groups
 - adding resource groups to True Copy, 287–289
 - configured MySQL, 305–306
 - creating Oracle Data Guard, 294–298
 - creating/starting Geographic Edition, 260–263
 - creating TrueCopy, 284–286
 - null data replication type, 312
 - placing resource groups in, 277–280
 - properties of SBP, 308–310
 - replication components and, 264–265
 - resource group configuration and, 268–272
 - SRDF replication resources and, 280–281
 - starting/stopping, 313–315

- switchover performed on, 316
- takeover performed on, 317–319
- Protocols, cluster private network, 103
- Provider, of HA services, 39
- Proxy file system (PxFS)
 - accessing file systems, 68
 - architecture, 69–70
 - cache coherence, 72
 - campus cluster and, 77
 - compared to NFS, 77–79
 - at file system level, 70–71
 - files/HA service, 71–72
 - HA invocation and, 41–42
 - I/O support, 72–73
 - Mount subsystem, 74
 - Oracle RAC files in, 68
 - performance, 76–77
- pts_partnership-name format, 253
- Public adapter utilization telemetry, 187
- Public network monitoring daemon, 23
- Public networking
 - cluster, 86–88
 - configuring IPMP for, 22–23
 - jumbo frames/link aggregation, 24–26
 - monitoring daemon, 23
 - overview, 21–22
 - using VLAN tagging, 24
 - zone clusters supporting, 160–162
- “Pull” model, MySQL, 300
- Pure service, 89–90
- PxFS. *See* Proxy file system (PxFS)

Q

- qd_userd daemon, 127
- Quorum devices
 - disk devices, 46
 - general behavior, 47–48
 - overview, 46
 - recommendations, 54
 - reservation mechanisms, 48–49
- Quorum monitoring, 53
- Quorum server, 46, 50–51
- Quorum subsystem, 37

R

- R1_SID property, 281
- R2_SID property, 281

RAID devices

- data protection and, 13–14
- disaster recovery and, 228
- of Oracle Solaris ZFS storage pool, 17
- zone clusters supporting access to, 158–159
- Random/sequential write performance, 32–33
- Raw-disk devices, 313
- rawdsk-type device groups
 - matching SRDF device groups, 277
 - multiple DID devices into, 64
- RBAC. *See* Role-based access control (RBAC)
- RDBMS (relational database management systems), 291–293
- Read-only restriction, 14, 15, 20, 43, 59, 145
- Recommendations, quorum device, 54
- Reconfiguration process, cluster, 55–56, 57
- Reconfiguration subsystem, 37
- Recovery
 - after GIN node failure, 98–99
 - options, for Oracle 11g database, 366
 - of Sun QFS file systems, 80
- Recovery Point Objective (RPO)
 - disaster recovery and, 223, 225, 227
 - storage replication and, 18
- Recovery Time Objective (RTO), 223, 225, 227
- Redundancy, data, 13
- Redundant components, 2
- Registration keys, 48
- Relational database management systems (RDBMS), 291–293
- Reliability, of Solaris Cluster software, 1–5
- remote_database_name property, 296
- remote_db_service_name property, 296
- remote_logical_host property, 272
- remote_rac_proxy_svr_rg_name property, 295
- remote_service_password replication component property, 311
- remove_app_rg_args protection group property, 309
- remove_app_rg_script protection group property, 310
- remove_config_args replication component property, 311
- remove_config_script protection group property, 310

- replctl command, 40
- Replica Manager, 39, 56
- Replication components
 - EMC Symmetrix Remote Data Facility.
 - See EMC Symmetrix Remote Data Facility (SRDF)
 - Hitachi replication software. *See* Hitachi Data Systems Universal Replicator; Hitachi TrueCopy
 - MySQL replication. *See* MySQL replication
 - null (none) data replication type, 312
 - Oracle Data Guard. *See* Oracle Data Guard software
 - overview, 263–265
 - StorageTek AVS software, 265–272
- Replication, for data protection
 - application-based, 232–233
 - file system protection, 233–235
 - host-based, 230–231
 - overview, 227–228
 - storage-based, 228–229
- Replication modules
 - additional technology for, 329–330
 - debugging information for, 327–328
 - implementation, 252
- Replication resources
 - AVS monitoring, 272
 - Hitachi replication software monitoring, 289–291
 - MySQL monitoring, 306
 - Oracle Data Guard monitoring, 299
 - replication components and, 264
 - SBP-type, 312
 - SRDF monitoring, 280–281
- Replication status resource group, 308
- replication_mode property, 296
- replication_role property
 - of SRDF replication resource, 280
 - of TrueCopy replication resource, 290
- Repositories, for cluster configuration, 42–44
- Required Java SE Software Development Kit, 340
- Required shared component upgrades, 341
- Resilience, of Solaris Cluster software, 1–5
- Resilience, Solaris Cluster software, 103
- Resource group affinities
 - in different zones, 122–123
 - setting, 119–122
 - types of, 118–119
- Resource Group Manager (RGM) subsystem
 - function of, 37–38
 - in Geographic Edition, 248
 - messages on state changes, 204–205
 - tracking resource utilization, 187
- Resource groups
 - adding to TrueCopy protection groups, 287–289
 - AVS configuration, 268–272
 - binding to processor pool project, 184–185
 - containing logical host resources, 86–87
 - dependency/affinity options, 118–123
 - for Geographic Edition software, 248–249
 - MySQL configuration, 305–306, 307
 - Oracle Data Guard configuration, 299
 - overview, 117
 - property settings of, 188–189
 - SRDF configuration, 277–280
 - TrueCopy configuration, 289, 290
- Resource Management API (RMAPI)
 - application integration and, 202
 - creating new resource types, 218
- Resource type registration (RTR) file
 - for creating GDS subclass, 208–212
 - SUNW.LogicalHostname, 216–217
- Resource types
 - callback methods, 109
 - function of, 108–109
 - instantiations of. *See* Resources with shared address resource, 110–111
- Resource types, creating
 - application suitability, 201–202
 - creating GDS subclass for, 208–212
 - Generic Data Service and, 203–207
 - integrating application-specific logic, 214–216
 - overview, 201
 - registration file, 216–217
 - Resource Management API/DSDL, 218
 - sdcdbuilder GUI, 208–211
 - supporting new applications, 207
 - tuning/troubleshooting, 220–221
 - utilities for building custom data services, 219–220
- Resources
 - dependency properties of, 114–115

- failover/scalable/multi-master, 116
 - for Geographic Edition software, 248–249
 - overview, 112
 - per-node properties, 115–116
 - state changes, 113
- Restart dependency, of resources, 112
- Retry_count property, 205–206
- Retry_interval property, 205–206
- RG_affinities property, 118–123
- RG_dependencies property, 118–123
- RGM. *See* Resource Group Manager (RGM)
 - subsystem
- rgmd daemon, 107, 117, 128
- RG_slm_type property, 188
- RMAPI. *See* Resource Management API (RMAPI)
- Role-based access control (RBAC)
 - for cluster management, 178–180
 - Geographic Edition cluster management, 325
- role property, of AVS replication resource, 272
- Rolling failures
 - data consistency after, 19–20
 - ZFS file system and, 234–235
- Rolling upgrades, 191–192
- Root (/) file system
 - application data backup and, 200
 - creating snapshot, 199
 - of HA containers, 243
 - mirrored zpool for, 344
- Root directory
 - non-replication of, 244
 - for Oracle Solaris zones, 145
- Root disk
 - backup, 198–199
 - partition requirement, 164–165
- Root state replica databases, 165
- round_robin property
 - deterministic distribution with, 92–94
 - packet distribution and, 91–92
- rpc.fed daemon, 107, 128
- rpc.mdcommnd daemon, 128
- rpc.metaclld daemon, 128
- rpc.metad daemon, 128, 129
- rpc.metamedd daemon, 128
- rpc.pmfdd daemon, 128
- RPO. *See* Recovery Point Objective (RPO)
- RTO. *See* Recovery Time Objective (RTO)
- RTR. *See* Resource type registration (RTR)
 - file
- rtreg_proxy_serverd daemon, 128
- S**
- SAM-QFS software
 - high capacity/file system recovery, 80
 - metadata storage/vnode interface, 80–81
 - overview/volume management, 79
 - paged/direct I/O support, 79–80
 - shared file system support/benefits, 81–82
 - shared QFS file system, 82–84
 - supporting clients outside cluster, 84
- SAP enqueue resource group, 122
- SAP replica server resource group, 122
- SBP. *See* Script-Based Plug-In (SBP) module
- Scalable data service, 105
- Scalable device group/mount-point
 - resources, 125
- Scalable resources
 - features of, 116
 - generic data service and, 203
 - holding file system/web server resources, 381
 - packet distribution and, 91
- Scalable services
 - four-node cluster hosting, 90
 - IPsec SA failover for, 97–98
 - protecting with AVS software, 398–405
 - pure, 89–90
 - sticky, 90–91
- Scalable storage topology, 8–9
- Scalable web services. *See* Sun Java System
 - Web server service, creating
- scds_prefix, for DSDL function names, 218
- scdsbuilder GUI, 212–214
- scdsbuilder tool, 202
- scinstall command, 192–194
- scprivipd daemon, 128
- scqcmd daemon, 128–129
- Script-Based Plug-In (SBP) module
 - configuration file, 304
 - creating additional replication
 - technology, 329–330
 - limitations of, 308
 - for MySQL replication, 306–308
 - protection group properties, 308–310

- SCSI-2 Reservation protocol/Persistent Group Reservation emulation (SCSI-2 PGR), 48–49
- SCSI-3 Persistent Group Reservation (PGR) protocol, 48
- `scsi-initiator-id` settings, 7
- SCSI (Small Computer System Interface)
 - fencing, 50, 58
 - storage devices, 7
- SCTP (Stream Control Transmission Protocol), 94–97
- `scversions -c` command, 191–192
- `sc_zonesd` daemon, 129
- Secondary (disaster recovery) sites
 - application-based replication and, 232
 - combined host-based mirroring/replication and, 241–242
 - connection topologies and, 236–237
 - host-based replication and, 231
 - migration to, 226–227
 - mirror consistency in, 282
 - mixing storage-based/host-based replication and, 239–241
 - Oracle Data Guard configuration and, 294–295
 - protection groups at, 313, 315
 - storage-based replication and, 229, 238
- Secondary I/O paths, 60, 63
- Secondary IT infrastructure site, 224, 225
- Secondary providers, for HA services, 39–40
- Secure by Default configuration, 168–169
- Security architecture, zone-cluster, 156–157
- Security association (SA), 97–98
- Security isolation, 151–152
- Segmented files, 81
- Sequential/random write performance, 32–33
- Server configurations, 6–7
- Server log file locations, 385
- Server pull communication technique, 73
- Server-to-storage connectivity topologies
 - clustered pair, 9–10
 - N+1, 10–11
 - overview, 8
 - Pair+M, 11–12
 - scalable, 8–9
- Service domains, 137
- Service-level agreements (SLAs)
 - business-critical services and, 224
 - N+1 topology and, 10–11
- Service-level management features
 - gathering telemetry, 185–189
 - overview, 183–185
- Service outages
 - with dual-partition upgrades, 194
 - installing Oracle on both cluster nodes and, 358
 - during node upgrades, 196–198
 - with rolling upgrades, 191
- Services, migration of
 - automated vs. automatic, 226–227
 - switchover as, 315–316
 - takeover as, 317–319
- Settings
 - for Sun Java System administration server, 384
 - for Sun Java System Web Server CGI directory, 386
 - for Sun Java System web server instance, 385
- Severe data corruption, with split-brain condition, 44
- SGA (System global area), 123–124
- Shared addresses, 95
- Shared component upgrades, 341
- Shared IP networking options, 145
- Shared (multi-owner) disksets, 15
- Shared QFS file systems
 - benefits of, 82
 - configuring/mounting, 82–84
 - function of, 68, 81
 - implementing of, 82–84
 - overview, 81–82
 - protecting with replication, 233–234
- Shared storage, 7, 47, 191
- SID (Storage unit identifier), 274–277
- Simple Network Management Protocol (SNMP)
 - adapter trap port, 182
 - management information bases, 182–183
 - overview, 180
- Single-node cluster system, creating, 335–339
- Single-owner (nonshared) disksets, 14–17
- Single point of administration, 153
- Single point of failure, 6–7, 99
- Single-threaded/multi-threaded write performance, 32–33
- Single writer-multiple reader model, 72
- SLAs. *See* Service-level agreements (SLAs)

- Slave (backup) database, 300
- Small Computer System Interface (SCSI), 7, 50, 58
- smpatch command, 343
- Snapshots. *See also* AVS snapshot files
 - application-based replication and, 232
 - HA framework and, 39
 - programming interface and, 54
 - of root (/) file system, 198–199
 - of UFS file systems, 200
 - ZFS, 224
- snapshot_volume property, 272
- SNMP. *See* Simple Network Management Protocol (SNMP)
- Software infrastructure, Geographic Edition
 - Common Agent Container (CAC)
 - modules, 249–251
 - event propagation/CLI, GUI, module implementation, 252
 - overview, 248
 - resources/resource groups, 248–249
 - storage of configuration information, 252–253
- Software quorum protocol, 49–50
- Solaris Cluster CLI commands, 174
- Solaris Cluster Clusterware framework
 - resource, 125
- Solaris Cluster, environment
 - backing up clusters, 198–201
 - cluster management. *See* Cluster management
 - cluster monitoring, 180–183
 - creating new resource types. *See* Resource types, creating
 - installing/configuring software, 169–172
 - installing operating system, 163–166
 - overview, 163
 - patching/upgrading. *See* Patching/upgrading
 - securing operating system, 166–169
 - service-level management/telemetry, 183–189
 - time synchronization, 172–173
- Solaris Cluster, features
 - architecture, 35–38
 - cluster configuration repository (CCR), 42–44
 - cluster membership. *See* Cluster membership
 - cluster private interconnect, 99–105
 - data service/application agents. *See* Data services/application agents
 - devices. *See* Devices
 - file systems. *See* File systems
 - HA framework, 39–42
 - Object Request Broker subsystem, 38–39
 - overview, 35
 - reconfiguration, 55–56
 - storage device fencing subsystem, 57–59
 - Version Manager subsystem, 56
- Solaris Cluster Geographic Edition, features
 - creating trust between clusters, 253–254
 - heartbeat messages, 257–259
 - overview, 247
 - partnerships between clusters, 254–257
 - protecting Oracle RAC databases, 313
 - protection groups, 260–263
 - replication components. *See* Replication components
 - required software infrastructure, 248–253
 - starting/stopping protection groups, 313–315
 - switchover/takeover, 315–319
- Solaris Cluster Geographic Edition implementations
 - configuration process, 395–398
 - overview, 395
 - protecting scalable web services, 398–405
- Solaris Cluster Geographic Edition, managing
 - cluster management, 324–326
 - creating data replication modules, 329–330
 - installing/removing, 321–324
 - overview, 321
 - patching/upgrading, 324
 - troubleshooting, 327–329
- Solaris Cluster Geographic Edition, overview
 - architecture, 226–227
 - choosing disaster recovery solution, 224–225
 - connection topologies, 235–237
 - data protection with replication. *See* Replication, for data protection
 - introduction, 223
 - reasons for disaster recovery solution, 223–224

- Solaris Cluster Geographic Edition,
 - overview (*continued*)
 - third-party disaster recovery framework, 225
 - three-data-center architecture, 237–242
 - virtualization technologies and, 242–245
- Solaris Cluster implementations
 - configuring software, on first node, 344–352
 - configuring software, on second node, 352–356
 - creating HA-Oracle 11g Release 1 database. *See* Oracle 11g Release 1 database, creating
 - creating HA-Oracle database instance, 387–393
 - creating scalable web services. *See* Sun Java System Web server service, creating
 - creating single-node cluster system, 335–339
 - installing Solaris 10 OS, 331–335
 - installing Solaris Cluster software, 339–342
 - overview, 331
 - setting up Solaris Cluster telemetry, 372–377, 378
- Solaris Cluster Manager
 - function of, 38
 - GUI, 175–177
 - wizards, 177–178
- Solaris Cluster Open Storage Program, 13
- Solaris Cluster, overview
 - architecture, 6–12
 - campus/metro clusters, 28–29
 - introduction, 1
 - private networking, 26–28
 - protecting data. *See* Data protection
 - public networking, 21–26
 - separating data centers, 30–34
 - value of, 1–5
- Solaris Cluster RAC/storage proxy resources, 126
- Solaris Cluster scalable device group resources, 125
- Solaris Cluster scalable mount-point resources, 125
- Solaris Cluster telemetry
 - displaying data collected, 378
 - setting up, 372–377
- Solaris Flash archive file, 163
- Solaris hosts
 - cluster nodes and, 132–133
 - zone clusters and, 150–151
- Solaris JumpStart, 163
- Solaris Live Upgrade
 - dual-partition upgrades and, 194
 - features of, 164
 - standard/rolling upgrades and, 191–192
 - updating with, 194–195
 - upgrading nodes and, 196–198
- Solaris operating system (OS)
 - hardening, 168–169
 - installation, 163–164, 331–335
 - minimization, 166–167
 - overview, 6
 - planning for upgrades, 165–166
 - public networking and, 23
 - root disk partition requirement, 164–165
 - securing communications, 169
- Solaris operating system (OS), securing
 - OS hardening, 168–169
 - OS minimization, 166–167
 - overview, 166
 - securing network communications, 169
- Solaris Resource Manager, 152–153
- Solaris root user, 252–253
- Solaris Security Toolkit, 168
- Solaris software groups, for installation, 167
- Solaris Volume Manager
 - data protection and, 13
 - disksets with, 14–15
 - root state replica databases, 165
 - state replica majority of, 15–17, 53
- SPARC architecture
 - Oracle VM server for. *See* Oracle VM server for SPARC
 - for quorum server, 50
 - for server configurations, 6
- Spare I/O paths, 60
- Spare providers, for HA services, 39–40
- Sparse-root zone, 145
- Split-brain conditions
 - in cluster membership, 44

- features of, 45
 - private interconnect failure and, 103
 - SCSI-2 PGR and, 49
 - software quorum and, 49–50
- SRDF. *See* EMC Symmetrix Remote Data Facility (SRDF)
- SRDF_group property, 280
- Standard Solaris /dev/rdisk devices, 66
- Standard upgrades, 190–191
- Standby databases
- Oracle Data Guard and, 291–292
 - read-only restriction on, 232
- standby_type property, 295
- "Star" topology, 100
- start and stop methods
- controlling service migration, 248
 - disk path monitoring and, 65
 - HA containers and, 147
 - resource-type, 115, 221, 226
- start_command property
- function with new program, 214–216
 - of GDS agent, 203–204
 - of GDS subclass, 208–211
- start_replication_args replication component property, 311
- State reconstruction, 71–72
- State replica database, 14
- State replica majority, 15–17
- STATUS fence level, 282
- stop_command property
- function with new program, 214–216
 - of GDS agent, 203
 - of GDS subclass, 208–211
- stop_replication_args replication component property, 311
- stop_replication_script protection group property, 310
- Storage-based replication, for data protection
- alone, in 3DC configuration, 238
 - in cluster membership, 47
 - disaster recovery and, 224
 - EMC Symmetrix range/Hitachi Universal Replicator, 228–229
 - file system protection with, 233–235
 - with host-based replication, 239–241
 - overview, 18–19
 - parallel application restrictions, 20
 - protecting Oracle RAC databases, 313
 - reasons for, 20–21
 - rolling failures and, 19–20
 - with Volume Managers, 229
- Storage connectivity, 7–8, 13
- Storage device fencing subsystem
- disk failfast driver, 58
 - disk fencing, 57–58
 - NAS device fencing, 59
 - overview, 57
- Storage devices, zone clusters and, 158–159
- Storage Fencing subsystem, 37
- Storage mechanism, for Oracle 11g database, 365
- Storage resources, 154
- Storage unit identifier (SID), 274–277
- StorageTek Availability Suite (AVS)
- communications failure and, 235
 - configuring AVS replication, 401–405
 - for host-based replication, 230–231
 - initialization files, 266–267
 - installing software, 399–400
 - monitoring, 272
 - protection groups and, 262–263
 - replication components and, 265
 - resource group configuration, 268–272
- st_pg_protection-group-name format, 253
- Stream Control Transmission Protocol (SCTP), 94–97
- Strong dependency, of resources, 112
- Strong negative resource group affinity, 119, 121, 122
- Strong positive resource group affinity/with failover delegation, 119–120
- subclass, of GDS, 208–212
- Subnets
- choosing allocation, 104–105
 - configuring NICs to multiple, 24
 - public networking and, 21–22
- Sun Cluster Manager (now Solaris Cluster Manager), 1, 38
- Sun Fire V240 servers, 344–352
- Sun Java System Web server service, creating
- completing configuration, 386–387
 - configuring settings, 384–386
 - creating zones for web servers, 378–382

- Sun Java System Web server service,
 - creating (*continued*)
 - installing software, 382–384
 - overview, 377
 - Sun Logical Domains. *See also* Oracle VM server for SPARC, 244–245
 - Sun Management Center software, 181, 183
 - Sun QFS file systems
 - high capacity/fast recovery, 80
 - metadata storage/vnode interface, 80–81
 - paged/direct I/O support, 79–80
 - shared file system support/benefits, 81–82
 - shared QFS file system on Solaris Cluster system, 82–84
 - supporting clients outside cluster, 84
 - volume management, 79
 - Sun StorEdge Traffic Manager, 7, 13
 - Sun Trunking software, 25
 - SUNWCuser package group, 167
 - SUNW.gds resource
 - function of, 112
 - for Oracle Enterprise Manager database control service, 358–360, 367, 370
 - registering, 203
 - SUNW.HAStoragePlus resource
 - function of, 111–112
 - in Oracle Solaris ZFS file system, 85
 - SUNW.ldom resource
 - guest domain failure and, 142
 - RTR file for, 208–212
 - SUNW.LogicalHostname resource
 - function of, 110
 - RTR file, extract from, 216–217
 - SUNW.SharedAddress resource
 - function of, 110
 - GIN node failure and, 98
 - SUNWsprot package, 358
 - Switch-connected private interconnects, 27–28
 - switchover_args replication component
 - property, 311
 - Switchovers, 315–316
 - switchover_script protection group
 - property, 310
 - Symmetrix command-line (SYMCLI)
 - programs, 274
 - Symmetrix Remote Data Facility (SRDF). *See* EMC Symmetrix Remote Data Facility (SRDF)
 - Synchronization, time, 172–173
 - Synchronous replication
 - mirror consistency in, 282
 - with storage-based replication, 228
 - syncsa_serverd daemon, 129
 - sysdba_password property, 294–298
 - sysdba_username property, 296
 - syslogd output, 53
 - System global area (SGA), 123–124
- T**
- takeover_args replication component
 - property, 311
 - Takeovers
 - defined, 315
 - on protection groups, 317–319
 - takeover_script protection group
 - property, 310
 - TCP/IP (Internet Protocol Suite), 103–104
 - Telemetry
 - displaying data collected, 378
 - service-level management feature
 - gathering, 185–189
 - setting up on Solaris Cluster system, 372–377
 - Third-party disaster recovery framework, 225
 - Third-site mediator, 16–17
 - Three-data-center (3DC) configuration
 - combining host-based mirroring/replication, 241–242
 - mixing storage-based/host-based replication, 239–241
 - overview, 237–238
 - storage-based replication only, 238
 - Time synchronization, across cluster nodes, 172–173
 - Timeout period, 51
 - Topology(ies)
 - connection, 235–237
 - private interconnect, 99–100
 - Traffic, private interconnect, 100–101
 - transaction_state data structure, 40
 - Transport subsystem
 - function of, 37
 - reconfiguration process and, 55
 - Troubleshooting
 - debugging information for, 327–329
 - unhealthy applications, 220–221

TrueCopy. *See* Hitachi TrueCopy
 Trunking, 25
 Trust, between clusters, 253–254, 255–257, 395, 397
 Tuning, unhealthy applications, 220–221
 Two-data-center configuration
 host-based replication in, 231
 storage-based replication in, 228–229

U

UDP (User Datagram Protocol)-based
 service, 87–88
 udp_session_timeout property, 93–94
 Uneven cluster partitions, 52–53
 Uninstalling, Geographic Edition, 324
 Universal Installer Summary, 362
 Universal Replicator. *See* Hitachi Data
 Systems Universal Replicator
 UNIX mount command, 74
 Unmanaged state, of resource groups, 117
 updatemanager program, 343
 Upgrade methods
 dual-partition, 192–194
 factors in choosing, 194
 overview/checklist, 190
 rolling, 191–192
 Solaris Live Upgrade, 194–195
 standard, 190–191
 Upgrades. *See also* Patching/upgrading
 of Geographic Edition software, 324
 planning for, 165–166
 User Datagram Protocol (UDP)-based
 service, 87–88
 /usr/cluster/bin directory, 173

V

Validate_command property
 function with new program, 214–216
 GDS and, 203, 208–209
 Value, of Solaris Cluster software
 cluster framework benefits, 4–5
 HA/HPC/fault-tolerant systems, 3
 HA vs. disaster recovery, 4
 overview, 1–3
 Verification of system requirements, 341
 Veritas Volume Manager, 13–15
 Version Manager subsystem, 37, 56

Virtual hosts. *See* Logical (virtual) hosts
 Virtual nodes, 155–156
 Virtualization technologies, combining
 "black-box" vs. fine-grained control, 133–
 134
 defining cluster, 133
 defining cluster node, 132–133
 dynamic system domains, 134–136
 Oracle Solaris zones, 143–149
 Oracle VM server for SPARC. *See* Oracle
 VM server for SPARC
 overview, 131
 zone clusters. *See* Zone clusters
 Virtualization technologies, Geographic
 Edition
 with logical domains, 244–245
 with Oracle Solaris Zones, 243–244
 overview/dynamic system domains, 242
 Virtualized environments, 232–233
 Virtualized I/O, 140–141
 VLAN (Virtual LAN) tagging, 24
 vnode/VFS interface, 80–81
 Volume Managers
 storage-based replication with, 228
 in Sun QFS file systems, 79
 zone clusters and, 159
 Votes, cluster node, 45–46

W

Weak dependency, of resources, 112
 Weak negative resource group affinity, 119,
 120–121
 Weak positive resource group affinity, 118,
 120, 122
 Web server instance, Sun Java System, 383,
 385
 Web service, Geographic Edition AVS
 configuration for, 268–272
 Whole-root zone, 145
 Wizards, Solaris Cluster Manager, 177–178

X

X11 program, 203–204
 x64 architecture, 6, 50
 xeyes service, 203–207
 xntpd daemon, 172

- Z**
- ZFS file system. *See* Oracle Solaris ZFS file system
 - Zone-cluster membership, 55
 - Zone Cluster Membership Monitor (ZCMM), 155–156
 - Zone clusters
 - administrative work reduction/storage resources, 154
 - application fault isolation, 152
 - architecture, 154–158
 - creating, 388–393
 - dedicated cluster model/single administration point, 153
 - features of, 133, 150–151
 - networks, 160–162
 - resource management, 152–153
 - security isolation, 151–152
 - storage devices, 158–159
 - zonecfg command, 143
 - Zones
 - cluster nodes as, 132
 - creating for web servers, 378–382
 - Oracle Solaris. *See* Oracle Solaris zones (containers)
 - setting affinities between, 122–123
 - Zpools
 - constructing from DID devices, 63
 - upgrading nodes and, 196–197
 - ZFS file system and, 17, 85, 234–235