



FailSafe™ Programmer's Guide for
SGI® InfiniteStorage

007-3900-007

CONTRIBUTORS

Written by Lori Johnson

Illustrated by Chrystie Danzer, Dany Galgani, and Chris Wengelski

Engineering contributions by Scott Henry, Vidula Iyer, Herb Lewis, Michael Nishimoto, Kevan Rehm, Hugh Shannon Jr., Bill Sparks, Paddy Sreenivasan, Dan Stekloff, Rebecca Underwood, and Manish Verma

COPYRIGHT

© 1999-2003 Silicon Graphics, Inc. All rights reserved; provided portions may be copyright in third parties, as indicated elsewhere herein. No permission is granted to copy, distribute, or create derivative works from the contents of this electronic documentation in any manner, in whole or in part, without the prior written permission of Silicon Graphics, Inc.

LIMITED RIGHTS LEGEND

The electronic (software) version of this document was developed at private expense; if acquired under an agreement with the USA government or any contractor thereto, it is acquired as "commercial computer software" subject to the provisions of its applicable license agreement, as specified in (a) 48 CFR 12.212 of the FAR; or, if acquired for Department of Defense units, (b) 48 CFR 227-7202 of the DoD FAR Supplement; or sections succeeding thereto. Contractor/manufacturer is Silicon Graphics, Inc., 1600 Amphitheatre Pkwy 2E, Mountain View, CA 94043-1351.

TRADEMARKS AND ATTRIBUTIONS

Silicon Graphics, SGI, the SGI logo, IRIS, IRIX, and XFS are registered trademarks and CXFS, FailSafe, IRIS FailSafe, SGI FailSafe, and SGI Linux, are trademarks of Silicon Graphics, Inc., in the United States and/or other countries worldwide.

Informix is a registered trademark of IBM. Linux is a registered trademark of Linus Torvalds, used with permission by Silicon Graphics, Inc. Netscape is a trademark of Netscape Communications Corporation. Oracle is a registered trademark of Oracle Corporation. Sybase is a trademark of Sybase, Inc.

New Features in This Guide

This revision adds the following:

- A list of `ha_exec2` exit codes; see "Implementing Timeouts and Retrying a Command" on page 12.
- Information about upgrading from FailSafe 1.2 has been moved to the *Migrating from IRIS FailSafe 1.2 to IRIS FailSafe 2.1.X* white paper.

Record of Revision

Version	Description
002	December 1999 Published in conjunction with the latest IRIS FailSafe 2.0 rollup patch. It supports IRIX 6.5.9 and later.
003	October 2000 Supports IRIS FailSafe 2.1.
004	April 2001 Supports IRIS FailSafe 2.1.1 release and IRIX 6.5.12 or later.
005	November 2001 Supports IRIS FailSafe 2.1.2 and IRIX 6.5.14 or later.
006	April 2002 Supports IRIS FailSafe 2.1.3 and IRIX 6.5.16 or later.
007	November 2003 Supports IRIS FailSafe 2.1.6 and IRIX 6.5.22 or later.

Contents

About This Guide	xvii
Audience	xvii
Related Documentation	xvii
Conventions Used in This Guide	xix
Obtaining Publications	xix
Reader Comments	xx
1. Introduction	1
Plug-ins	1
Characteristics that Permit an Application to be Highly Available	3
Overview of the Programming Steps	4
Administrative Commands for Use in Scripts	5
2. Writing the Action Scripts and Adding Monitoring Agents	7
Set of Action Scripts	7
Understanding the Execution of Action Scripts	8
When Action Scripts are Executed	9
Multiple Instances of a Script Executed at the Same Time	9
Differences between the exclusive and monitor Scripts	10
Successful Execution of Action Scripts	11
Failure of Action Scripts	12
Implementing Timeouts and Retrying a Command	12
Sending Signals	13
Preparation	14
007-3900-007	vii

Is Monitoring Necessary?	15
Types of Monitoring	16
What are the Symptoms of Monitoring Failure?	16
How Often Should Monitoring Occur?	16
Examples of Testing for Monitoring Failure	17
Script Format	18
Header Information	18
Set Local Variables	19
Read Resource Information	20
Exit Status	20
Basic Action	21
Set Global Variables	22
Verify Arguments	22
Read Input File	22
Complete the Action	23
Steps in Writing a Script	23
Examples of Action Scripts	24
start Script	24
stop Script	26
monitor Script	28
exclusive Script	31
restart Script	33
Monitoring Agents	34
3. Creating a Failover Policy	37
Contents of a Failover Policy	37
Failover Domain	37
Failover Attributes	39

Failover Scripts	41
ordered	41
round-robin	44
Creating a New Failover Script	48
Failover Script Interface	48
Creating a Failover Policy that Returns the Resource Group to the Same Node	49
Example Failover Policies	49
N+1 Configuration	50
N+2 Configuration	51
N+M Configuration	53
4. Defining a New Resource Type	55
Information You Must Gather	56
Copying an Existing Resource Type to Create a New One	59
Creating a New Resource Type from Scratch	60
Using the FailSafe Manager GUI	61
Define a New Resource Type	61
Define Dependencies	65
Using <code>cmgr</code> Interactively	66
Using <code>cmgr</code> With a Script	71
Server-side Properties File	73
Property Formats	74
Example Properties File	74
Testing a New Resource Type	76
5. Testing Scripts	79
General Testing and Debugging Techniques	79
Debugging Notes	80

Testing an Action Script	81
Special Testing Considerations for the monitor Script	83
6. Example: Requiring Confirmation Before Failover	85
Appendix A. Starting the FailSafe Manager GUI	87
Launch Methods	87
Logging In	89
Making Changes from One Node	90
Appendix B. Using the Script Library	91
File Formats	91
Set Global Definitions	93
Global Variable	93
HA_HOSTNAME	93
Command Location Variables	93
HA_CMDSPATH	93
HA_PRIVCMDSPATH	93
HA_LOGCMD	93
HA_RESOURCEQUERYCMD	94
HA_SCRIPTTMPDIR	94
Database Location Variables	94
HA_CDB	94
Script Log Level Variables	94
HA_NORMLVL	94
HA_DBGLVL	94
Script Log Variables	95
HA_SCRIPTGROUP	95
HA_SCRIPTSUBSYS	95

Script Logging Command Variables	95
HA_LOGQUERY_OUTPUT	95
HA_DBGLOG	95
HA_CURRENT_LOGLEVEL	95
HA_LOG	96
Script Error Value Variables	96
HA_SUCCESS	96
HA_NOT_RUNNING	96
HA_INVALID_ARGS	96
HA_CMD_FAILED	96
HA_RUNNING	96
HA_NOTSUPPORTED	97
HA_NOCFGINFO	97
Check Arguments	97
Read an Input File	98
Execute a Command	99
Write Status for a Resource	100
Get the Value for a Field	101
Get the Value for Multiple Fields	101
Get Resource Information	102
Print Exclusivity Check Messages	105
Glossary	107
Index	119

Figures

Figure 2-1	Monitoring Process	34
Figure 3-1	<i>N</i> +1 Configuration Concept	50
Figure 3-2	<i>N</i> +2 Configuration Concept	52
Figure 3-3	<i>N</i> + <i>M</i> Configuration Concept	53
Figure 4-1	Specify the Name of the New Resource Type	62
Figure 4-2	Specify Settings for Required Actions	63
Figure 4-3	Change Settings for Optional Actions	64
Figure 4-4	Set Type-specific Attributes	65
Figure 4-5	Add Dependencies	66

Tables

Table 1-1	Provided Plug-Ins	2
Table 1-2	Optional Plug-Ins	2
Table 1-3	FailSafe Administrative Commands for Use in Scripts	6
Table 2-1	Execution of Action Scripts	9
Table 2-2	Differences Between the monitor and exclusive Action Scripts	11
Table 2-3	Successful Action Script Results	11
Table 2-4	Failure of an Action Script	12
Table 3-1	Failover Attributes	40
Table 4-1	Order Ranges	57
Table 4-2	Resource Type Order Numbers	57
Table A-1	GUI Platforms	89

About This Guide

This guide explains how to write your own *plug-in*, the set of scripts that are required to turn an application into a highly available service in conjunction with IRIS FailSafe 2.1.6 software.

This guide assumes that the FailSafe system has been configured as described in the *FailSafe Administrator's Guide for SGI InfiniteStorage*.

This guide supports IRIX 6.5.22 and later.

Audience

This guide is written for system programmers who are writing their own plug-ins for the FailSafe system. This software allows the failover of applications that are not handled by the base and optional plug-ins. Readers must be familiar with the operation and administration of nodes running FailSafe, with the applications that are to be failed over, and with the *FailSafe Administrator's Guide for SGI InfiniteStorage*.

Related Documentation

The following documentation is of interest:

- *FailSafe Administrator's Guide for SGI InfiniteStorage*
- *CXFS Administration Guide for SGI Infinite Storage*
- *Migrating from IRIS FailSafe 1.2 to IRIS FailSafe 2.1.X*

The man pages are as follows:

- `cdbBackup`
- `cdbRestore`
- `cmgr`
- `crsd`
- `failsafe`

- fs2d
- ha_cilog
- ha_cmds
- ha_cxfs
- ha_exec2
- ha_fsd
- ha_gcd
- ha_ifd
- ha_ifdadmin
- ha_macconfig2
- ha_srmd
- ha_statd2
- haStatus

Release notes are included with each FailSafe product. The names of the release notes are as follows:

Release Note	Product
cluster_admin	Cluster administration services
cluster_control	Cluster node control services
cluster_services	Cluster services
failsafe2	IRIS 2.x FailSafe
failsafe2_informix	FailSafe Informix
failsafe2_nfs	FailSafe NFS
failsafe2_oracle	FailSafe Oracle
failsafe2_samba	FailSafe Samba
failsafe2_web	FailSafe Netscape web

Conventions Used in This Guide

These type conventions and symbols are used in this guide:

Bold	Function names literal command-line arguments (options/flags)
Bold fixed-width type	Commands and text that you are to type literally in response to shell and command prompts, or highlighting of differences between releases
<i>Italics</i>	New terms, manual/book titles, commands, variable command-line arguments, filenames, and variables to be supplied by the user in examples, code, and syntax statements
Fixed-width type	Code examples, error messages, prompts, and screen text
#	IRIX shell prompt for the superuser (root)

Obtaining Publications

You can obtain SGI documentation as follows:

- See the SGI Technical Publications Library at <http://docs.sgi.com>. Various formats are available. This library contains the most recent and most comprehensive set of online books, release notes, man pages, and other information.
- If it is installed on your SGI system, you can use InfoSearch, an online tool that provides a more limited set of online books, release notes, and man pages. With an IRIX system, enter `infosearch` at a command line or select **Help > InfoSearch** from the Toolchest.
- On IRIX systems, you can view release notes by entering either `grelnotes` or `relnotes` at a command line.
- On Linux systems, you can view release notes on your system by accessing the `README.txt` file for the product. This is usually located in the `/usr/share/doc/productname` directory, although file locations may vary.
- You can view man pages by typing `man title` at a command line.

Reader Comments

If you have comments about the technical accuracy, content, or organization of this publication, contact SGI. Be sure to include the title and document number of the publication with your comments. (Online, the document number is located in the front matter of the publication. In printed publications, the document number is located at the bottom of each page.)

You can contact SGI in any of the following ways:

- Send e-mail to the following address:

techpubs@sgi.com

- Use the Feedback option on the Technical Publications Library Web page:

<http://docs.sgi.com>

- Contact your customer service representative and ask that an incident be filed in the SGI incident tracking system.

- Send mail to the following address:

Technical Publications
SGI
1600 Amphitheatre Parkway, M/S 535
Mountain View, California 94043-1351

SGI values your comments and will respond to them promptly.

Introduction

SGI FailSafe provides highly available services for as many as eight nodes in a cluster. These services are monitored by the FailSafe software. You can create additional services that are highly available by using the instructions in this guide to write your own *plug-in*, the set of scripts that are required to turn an application into a highly available service in conjunction with FailSafe software.

This chapter contains the following:

- "Plug-ins"
- "Characteristics that Permit an Application to be Highly Available" on page 3
- "Overview of the Programming Steps" on page 4
- "Administrative Commands for Use in Scripts" on page 5

For information about FailSafe terminology and the components, software layers, communication paths, and order of execution of action and failover scripts, see *FailSafe Administrator's Guide for SGI InfiniteStorage*.

Plug-ins

A *plug-in* is the set of software required to make an application highly available, including a resource type and action scripts. There are plug-ins provided with the base FailSafe release, optional plug-ins available for purchase from SGI, and customized plug-ins you can write using the instructions in this guide. See "Plug-ins" on page 1.

The following tables show the provided and optional FailSafe plug-ins and their associated resource types.

Table 1-1 Provided Plug-Ins

Provided Plug-In	Resource Type
CXFS file system	CXFS
IP addresses	IP_address
MAC addresses	MAC_address
XFS filesystems	filesystem
XLV logical volumes	volume
XVM volume manager	XVM

Table 1-2 Optional Plug-Ins

Optional Plug-In	Resource Type
IRIS FailSafe for DMF	DMF
IRIS FailSafe for NFS	NFS and statd_unlimited
IRIS FailSafe for Informix	INFORMIX_DB
IRIS FailSafe for Oracle	Oracle_DB
IRIS FailSafe for Samba	Samba
IRIS FailSafe for TMF	TMF
IRIS FailSafe for Web (Netscape)	Netscape_web

See the release notes for information about the specific releases that are supported.

If you want to create your own plug-in, or change the functionality of the provided failover scripts and action scripts by writing new scripts, you will use the instructions in this guide.

Note: If you require a customized plug-in but do not want to write it yourself, you can establish a contract with the Silicon Graphics Professional Services group to create customized scripts. See: <http://www.sgi.com/services/index.html>.

Characteristics that Permit an Application to be Highly Available

The characteristics of an application that can be made highly available are as follows:

- The application can be easily restarted and monitored. It should be able to recover from failures as does most client/server software. The failure could be a hardware failure, an operating system failure, or an application failure. If a node crashes and reboots, client/server software should be able to attach again automatically.
- The application must have a start and stop procedure. When the application fails over, the instances of the application are stopped on one node using the stop procedure and restarted on the other node using start procedure.

Avoid applications that are started as a daemon from `/etc/inetd.conf` because typically everything in `/etc/inetd.conf` is already running. Trying to automatically edit `/etc/inetd.conf` could cause errors for other daemons started by this file.

Many applications will have a start and stop procedure that belongs in the `/etc/init.d` directory. You can incorporate them into a custom `/var/ha/resources` script to appropriately start and stop the application. If the application also has a `chkconfig` flag, set it to `off`. The `chkconfig` flag should be set to `on` in the `/var/ha/resources` start script.

- The application does not depend on the `hostname` or any identifier that is specific to a node.
- The application can be moved from one node to another after failures.

If the resource has failed, it must still be possible to run the resource stop procedure. In addition, the resource must recover from the failed state when the resource start procedure is executed on another node.

Ensure that there is no affinity for a specific node.

- The application does not depend on knowing the primary host name (as returned by `hostname`); that is, those resources that can be configured to work with an IP address.
- Other resources on which the application depends can be made highly available. If they are not provided by FailSafe and its optional products (see "Plug-ins" on page 1), you must make these resources highly available, using the information in this guide.

Note: An application itself is not modified to make it highly available.

Overview of the Programming Steps

To make an application highly available, follow these steps:

1. Understand the application and determine the following:
 - The configuration required for the application, such as user names, permissions, data location (volumes), and so on. For more information about configuration, see the *FailSafe Administrator's Guide for SGI InfiniteStorage*.
 - The other resources on which the application depends. All interdependent resources must be part of the same resource group.
 - The resource type that best suits this application.
 - The number of instances of the resource type that will constitute the application. (Each instance of a given application, or resource type, is a separate resource.) For example, a web server may depend upon two filesystem resources.
 - The commands and arguments required to start, stop, and monitor this application (that is, the resources in the resource group).
 - The order in which all resources in the resource group must be started and stopped.
2. Determine whether existing action scripts can be reused. If they cannot, write a new set of action scripts, using existing scripts and the templates in `/var/cluster/ha/resource_types/template` as a guide. See Chapter 2, "Writing the Action Scripts and Adding Monitoring Agents" on page 7.
3. Determine whether the existing `ordered` or `round-robin` failover scripts can be reused for the resource group. If they cannot, write a new failover script. See Chapter 3, "Creating a Failover Policy" on page 37.
4. Determine whether an existing resource type can be reused. If none applies, create a new resource type or modify an existing resource. See Chapter 4, "Defining a New Resource Type" on page 55.

5. Configure the following in the cluster database (for more information, see the *FailSafe Administrator's Guide for SGI InfiniteStorage*):
 - Resource group
 - Resource type
 - Failover policy
6. Test the action scripts and failover script. See Chapter 5, "Testing Scripts" on page 79, and "Debugging Notes" on page 80.

Note: Do not modify the scripts included with the FailSafe release. New or customized scripts must have different names from the files included with the release.

Administrative Commands for Use in Scripts

Table 1-3 shows the administrative commands available with FailSafe for use in scripts.

Table 1-3 FailSafe Administrative Commands for Use in Scripts

Command	Purpose
ha_cilog	Logs messages to the <code>script_nodename</code> log files.
ha_execute_lock	Executes a command with a lock file that allows command execution to be serialized. The lock file prevents multiple instances of the same command from executing at the same time on a single node.
ha_exec2	Executes a command and retries the command on failure or timeout.
ha_filelock	Locks a file.
ha_fileunlock	Unlocks a file.
ha_ifdadmin	Communicates with the <code>ha_ifd</code> network interface agent daemon.
ha_http_ping2	Checks if a web server is running.
ha_macconfig2	Displays or modifies MAC addresses of a network interface.



Caution: Do not use the script in `/usr/sysadm/privbin`. These are internal commands that have a different command line parameter scheme. The functionality of these commands may change in the future. These commands are not documented.

Writing the Action Scripts and Adding Monitoring Agents

This chapter describes how to write the action scripts required for a plug-in and how to add monitoring agents. It discusses the following topics:

- "Set of Action Scripts"
- "Understanding the Execution of Action Scripts" on page 8
- "Preparation" on page 14
- "Script Format" on page 18
- "Steps in Writing a Script" on page 23
- "Examples of Action Scripts" on page 24
- "Monitoring Agents" on page 34

Set of Action Scripts

The *action scripts* are the set of scripts that determine how a resource is started, monitored, and stopped.



Caution: Multiple instances of scripts may be executed at the same time. For more information, see "Understanding the Execution of Action Scripts" on page 8.

The following set of action scripts can be provided for each resource type:

- `exclusive`, which verifies that the resource is not already running
- `start`, which starts the resource
- `stop`, which stops the resource
- `monitor`, which monitors the resource
- `restart`, which restarts the resource on the same node when a monitoring failure occurs

The `start`, `stop`, and `exclusive` scripts are required for every resource type.

Note: The `start` and `stop` scripts must be *idempotent*; that is, they have the appearance of being run once but can in fact be run multiple times. For example, if the `start` script is run for a resource that is already started, the script must not return an error.

A `monitor` script is required, but if you wish it may contain only a `return-success` function. A `restart` script is required if the application must have a restart ability on the same node in case of failure. However, the `restart` script may contain only a `return-success` function.

Understanding the Execution of Action Scripts

Before you can write a new action script, you must understand how action scripts are executed. This section covers the following topics:

- "When Action Scripts are Executed" on page 9
- "Multiple Instances of a Script Executed at the Same Time" on page 9
- "Differences between the `exclusive` and `monitor` Scripts" on page 10
- "Successful Execution of Action Scripts" on page 11
- "Failure of Action Scripts" on page 12
- "Implementing Timeouts and Retrying a Command" on page 12
- "Sending Signals" on page 13

When Action Scripts are Executed

Table 2-1 shows the circumstances under which action scripts are executed.

Table 2-1 Execution of Action Scripts

Script	Execution Conditions
<code>exclusive</code>	A resource group is made online by the user High-availability (HA) processes (<code>ha_cmsd</code> , <code>ha_gcd</code> , <code>ha_fsd</code> , <code>ha_srmd</code> , <code>ha_ifd</code>) are started
<code>start</code>	A resource group is made online by the user HA processes are started A resource group fails over
<code>stop</code>	A resource group is made offline HA processes are stopped A resource group fails over A node is shutdown or rebooted
<code>monitor</code>	A resource groups is online
<code>restart</code>	The <code>monitor</code> script fails

Multiple Instances of a Script Executed at the Same Time

Multiple instances of the same script may be executed at the same time. To avoid this problem, you can use the `ha_filelock` and `ha_execute_lock` commands to achieve sequential execution of commands in different instances of the same script.

For example, multiple instances of `xlv_assemble` should not be executed in a node at the same time. Therefore, the start script for volumes should execute `xlv_assemble` under the control of `ha_execute_lock` as follows:

```
${HA_CMDSPATH}/ha_execute_lock 30  
${HA_SCRIPTTMPDIR}/lock.volume_assemble \"/sbin/xlv_assemble -l  
-s${VOLUME_NAME} \"
```

All resources of the same resource type in a given resource group are passed as parameters to the action scripts.

The `ha_execute_lock` command takes the following arguments:

- Number of seconds before the command times out waiting for the file lock
- File to be used for locking
- Command to be executed

The `ha_execute_lock` command tries to obtain a lock on the file every second for *timeout* seconds. After obtaining a lock on the file, it executes the command argument. On command completion, it releases the lock on the file.

Differences between the `exclusive` and `monitor` Scripts

Although the same check can be used in `monitor` and `exclusive` action scripts, they are used for different purposes. Table 2-2 summarizes the differences between the scripts.

Table 2-2 Differences Between the monitor and exclusive Action Scripts

<code>exclusive</code>	<code>monitor</code>
Executed in all nodes in the cluster.	Executed only on the node where the resource group (which contains the resource) is online.
Executed before the resource is started in the cluster.	Executed when the resource is online in the cluster. (The <code>monitor</code> script could degrade the services provided by the HA server. Therefore, the check performed by the <code>monitor</code> script should be lightweight and less time consuming than the check performed by the <code>exclusive</code> script.)
Executed only once before the resource group is made online in the cluster.	Executed periodically.
Failure will result in resource group not becoming online in the cluster.	Failure will cause a resource group failover to another node or a restart of the resource in the local node.

Successful Execution of Action Scripts

Table 2-3 shows the state of a resource group after the successful execution of an action script for every resource within a resource group. To view the state of a resource group, use the FailSafe Manager graphical user interface (GUI) or the `cmgr` command.

Table 2-3 Successful Action Script Results

Event	Resource Group State	Script to Execute
Resource group is made online on a node	<code>online</code>	<code>start</code>
Resource group is made offline on a node	<code>offline</code>	<code>stop</code>
Online status of the resource group	(No effect)	<code>exclusive</code>
Normal monitoring of online resource group	<code>online</code>	<code>monitor</code>
Resource group monitoring failure	<code>online</code>	<code>restart</code>

Failure of Action Scripts

Table 2-4 shows the state of the resource group and the error state when an action script fails. (There are no `offline` states with errors.)

Table 2-4 Failure of an Action Script

Failing Script	Resource Group State	Error State
<code>exclusive</code>	<code>online</code>	<code>exclusivity</code>
<code>monitor</code>	<code>online</code>	<code>monitoring failure</code>
<code>restart</code>	<code>online</code>	<code>monitoring failure</code>
<code>start</code>	<code>online</code>	<code>srmd executable error</code>
<code>stop</code>	<code>online</code>	<code>srmd executable error</code>

When monitoring fails, FailSafe will stop monitoring. After recovering the resource, the system administrator must bring the resource group online again in order for FailSafe to resume monitoring it. For example, if a `start` script fails, the state of the resource and the resource group will be `online` and the error will be `srmd executable error`. FailSafe will attempt to move the resource group to other nodes in the application failover domain when the `start` script fails

Implementing Timeouts and Retrying a Command

You can use the `ha_exec2` command to execute action scripts using timeouts. This allows the action script to be completed within the specified time, and permits proper error messages to be logged on failure or timeout. The `retry` variable is especially useful in `monitor` and `exclusive` action scripts.

To retry a command, use the following syntax:

```
/usr/cluster/bin/ha_exec2 timeout_in_seconds number_of_retries command
```

For example:

```
${HA_CMDSPATH}/ha_exec2 30 2 "umount /fs"
```

The above `ha_exec2` command executes the `umount /fs` command line. If the command does not complete within 30 seconds, it kills the `umount` command and

retries the command. The `ha_exec2` command retries the `umount` command twice if it times out or fails.

The `ha_exec2` command executes the command string passed as a parameter. If the command string successfully completes execution and it returns an exit code, then that exit code is returned by `ha_exec2`. However, if there is a failure, the following special exit codes are returned by `ha_exec2`:

- 100: the command could not be executed
- 101: there was an invalid argument to `ha_exec2`
- 102: the `ha_exec2` command failed
- 103: the command timed out and was killed by `ha_exec2`
- 104: the command timed out and `ha_exec2` could not kill the command
- 105: the command exited with no error code

For more information, see the `ha_exec2` man page.

Sending Signals

You can use the `ha_exec2` command to send signals to specific process. A process is identified by its name or its arguments.

For example:

```
${HA_CMDSPATH}/ha_exec2 -s 0 -t "SYBASE_DBSERVER"
```

The above command sends signal 0 (checks if the process exists) to all processes whose name or arguments match the `SYBASE_DBSERVER` string. The command returns 0 if it is a success.

You should use the `ha_exec2` command to check for server processes in the `monitor` script instead of using the `ps -ef | grep` command line.

For more information, see the `ha_exec2` man page.

Preparation

Before you can write the action scripts, you must do the following:

- Understand the `scriptlib` functions described in Appendix B, "Using the Script Library" on page 91.
- Familiarize yourself with the script templates provided in the following directory:
`/var/cluster/ha/resource_types/template`
- Read the man pages for the following commands:
 - `cmgr`
 - `fs2d`
 - `ha_cilog`
 - `ha_cmsd`
 - `ha_exec2`
 - `ha_fsd`
 - `ha_gcd`
 - `ha_ifd`
 - `ha_ifdadmin`
 - `ha_macconfig2`
 - `ha_srmd`
 - `ha_statd2`
 - `haStatus`
- Familiarize yourself with the action scripts for other highly available services in `/var/cluster/ha/resource_types` that are similar to the scripts you wish to create.
- Understand how to do the following actions for your application:
 - Verify that the resource is running
 - Verify that the resource can be run

- Start the resource
- Stop the resource
- Check for the server processes
- Do a simple query as a client and understand the expected response
- Check for configuration file or directory existence (as needed)
- Determine whether or not monitoring is required (see "Is Monitoring Necessary?" on page 15). However, even if monitoring is not needed, a `monitor` script is still required; in this case, it can contain only a return-success function.
- Determine if a resource type must be added to the cluster database.
- Understand the vendor-supplied startup and shutdown procedures.
- Determine the configuration parameters for the application; these may be used in the action script and should be stored in the cluster database. Action scripts may read from the database.
- Determine whether the resource type can be restarted in the local node and whether this action makes sense.

Is Monitoring Necessary?

In the following situations, you may not need to perform application monitoring:

- Heartbeat monitoring is sufficient; that is, simply verifying that the node is alive (provided automatically by the base software) determines the health of the highly available service.
- There is no process or resource that can be monitored. For example, the SGI Gauntlet Internet Firewall software performs IP filtering on firewall nodes. Because the filtering is done in the kernel, there is no process or resource to monitor.
- A resource on which the application depends is already monitored. For example, monitoring some client-node resources might best be done by monitoring the file systems, volumes, and network interfaces they use. Because this is already done by the base software, additional monitoring is not required.



Caution: Beware that monitoring should be as lightweight as possible so that it does not affect system performance. Also, security issues may make monitoring difficult. If you are unable to provide a monitoring script with appropriate performance and security, consider a monitoring agent; see "Monitoring Agents" on page 34.

Types of Monitoring

There are two types of monitoring that may be accomplished in a `monitor` script:

- Is the resource present?
- Is the resource responding?

You can define multiple levels of monitoring within the `monitor` script, and the administrator can choose the desired level by configuring the resource definition in the cluster database. Ensure that the monitoring level chosen does not affect system performance. For more information, see the *FailSafe Administrator's Guide for SGI InfiniteStorage*.

What are the Symptoms of Monitoring Failure?

Possible symptoms of failure include the following:

- The resource returns an error code
- The resource returns the wrong result
- The resource does not return quickly enough

How Often Should Monitoring Occur?

You must determine the monitoring interval time and time-out value for the `monitor` script. The time-out must be long enough to guarantee that occasional anomalies do not cause false failovers. It will be useful for you to determine the peak load that the resource may need to sustain.

You must also determine if the `monitor` test should execute multiple times so that an application is not declared dead after a single failure. In general, testing more than once before declaring failure is a good idea.

Examples of Testing for Monitoring Failure

The test should be simple and complete quickly, whether it succeeds or fails. Some examples of tests are as follows:

- For a client/server resource that follows a well-defined protocol, the `monitor` script can make a simple request and verify that the proper response is received.
- For a web server application, the `monitor` script can request a home page, verify that the connection was made, and ignore the resulting home page.
- For a database, a simple request such as querying a table can be made.
- For NFS, more complicated end-to-end monitoring is required. The test might consist of mounting an exported file system, checking access to the file system with a `stat()` system call to the root of the file system, and undoing the mount.
- For a resource that writes to a log file, check that the size of the log file is increasing or use the `grep` command to check for a particular message.
- The following command can be used to determine quickly whether a process exists:

```
/sbin/killall -0 process_name
```

You can also use the `ha_exec2` command to check if a process is running.

The `ha_exec2` command differs from `killall` in that it performs a more exhaustive check on the process name as well as process arguments. `killall` searches for the process using the process name only. The command line is as follows:

```
/usr/cluster/bin/ha_exec2 -s 0 -t process_name
```

Note: Do not use the `ps` command to check on a particular process because its execution can be too slow.

Script Format

Templates for the action scripts are provided in the following directory:

```
/var/cluster/ha/resource_types/template
```

The template scripts have the same general format. Following is the type of information in the order in which it appears in the template scripts:

- Header information
- Set local variables
- Read resource information
- Exit status
- Perform the basic action of the script, which is the customized area you must provide
- Set global variables
- Verify arguments
- Read input file

Note: Action “scripts” can be of any form – such as Bourne shell script, Perl script, or C language program. The rest of this chapter discusses Korn shell.

The following sections show an example from the NFS `start` script.

Header Information

The header information contains comments about the resource type, script type, and resource configuration format. You must modify the code as needed.

Following is the header for the NFS start script:

```
#!/sbin/ksh

# *****
# *
# *          Copyright (C) 1998 Silicon Graphics, Inc.          *
# *
# *  These coded instructions, statements, and computer programs contain *
# *  unpublished proprietary information of Silicon Graphics, Inc., and *
# *  are protected by Federal copyright law. They may not be disclosed *
# *  to third parties or copied or duplicated in any form, in whole or *
# *  in part, without the prior written consent of Silicon Graphics, Inc. *
# *
# *  *****
#ident "$Revision: 1.24 $"

# Resource type: NFS
# Start script NFS

#
# Test resource configuration information is present in the database in
# the following format
#
# resource-type.NFS
#
```

Set Local Variables

The `set_local_variables()` section of the script defines all of the variables that are local to the script, such as temporary file names or database keys. All local variables should use the `LOCAL_` prefix. You must modify the code as needed.

Following is the `set_local_variables()` section from the NFS start script:

```
set_local_variables()
{
LOCAL_TEST_KEY=NFS
}
```

Read Resource Information

The `get_xxx_info()` function, such as `get_nfs_info()`, reads the resource information from the cluster database. `$1` is the test resource name. If the operation is successful, a value of 0 is returned; if the operation fails, 1 is returned.

The information is returned in the `HA_STRING` variable. For more information about `HA_STRING`, see Appendix B, "Using the Script Library" on page 91.

Following is the `get_nfs_info()` section from the NFS start script:

```
get_nfs_info ()
{
    ha_get_info ${LOCAL_TEST_KEY} $1
    if [ $? -ne 0 ]; then
        return 1;
    else
        return 0;
    fi
}
```

Call `ha_get_info` with a third argument of any value to obtain all attributes and dependency information for a resource from the cluster database. Use `ha_get_multi_fields` to retrieve specific dependency information. The resource dependency information is returned in the `$HA_FIELD_VALUE` variable.

Exit Status

In the `exit_script()` function, `$1` contains the `exit_status` value. If cleanup actions are required, such as the removal of temporary files that were created as part of the process, place them before the `exit` line.

Following is the `exit_script()` section from the NFS start script:

```
exit_script()
{
    ${HA_DBGLOG} "Exit: exit_script()";
    exit $1;
}
```

Note: If you call the `exit_script` function prior to normal termination, it should be preceded by the `ha_write_status_for_resource` function and you should use the same return code that is logged to the output file.

Basic Action

This area of the script is the portion you must customize. The templates provide a minimal framework.

Following is the framework for the basic action from the `start` template:

```
start_template()

# for all template resources passed as parameter
for TEMPLATE in $SHA_RES_NAMES
do
    #HA_CMD="command to start $TEMPLATE resource on the local machine" ;

    #ha_execute_cmd "string to describe the command being executed" ;

    ha_write_status_for_resource $TEMPLATE $SHA_SUCCESS;
done
}
```

Note: When testing the script, you will add the following line to this area to obtain debugging information:

```
set -x
```

For examples of this area, see "Examples of Action Scripts" on page 24.

Set Global Variables

The following lines set all of the global and local variables and store the resource names in `$HA_RES_NAMES`.

Following is the `set_global_variables()` function from the NFS start script:

```
set_global_variables()
{
    HA_DIR=/var/cluster/ha
    COMMON_LIB=${HA_DIR}/common_scripts/scriptlib

    # Execute the common library file
    . $COMMON_LIB

    ha_set_global_defs;
}
```

Verify Arguments

The `ha_check_args()` function verifies the arguments and stores them in the `$HA_INFILE` and `$HA_OUTFILE` variables. It returns 1 on error and 0 on success.

Following is the `ha_check_args()` section from the NFS start script:

```
ha_check_args $*;
if [ $? -ne 0 ]; then
    exit $HA_INVALID_ARGS;
fi
```

Read Input File

The `ha_read_infile()` function reads the input file and stores the resource names in the `$HA_RES_NAMES` variable. This function is defined in the `scriptlib` library. See "Read an Input File" on page 98.

Following is code from the NFS start script that calls the `ha_read_infile()` function:

```
# Read the input file and store the resource names in $HA_RES_NAMES
# variable

ha_read_infile;
```

Complete the Action

Each action script ends with the following, which performs the action and writes the output status to the `$HA_OUTFILE`:

```
action_resourcetype;

exit_script $HA_SUCCESS
```

Following is the completion from the NFS start script:

```
start_nfs;

exit_script $HA_SUCCESS;
```

Steps in Writing a Script



Caution: Multiple copies of actions scripts can execute at the same time. Therefore, all temporary filenames used by the scripts can be suffixed by `script.$$` in order to make them unique, or you can use the resource name because it must be unique to the cluster.

For each script, you must do the following:

- Get the required variables
- Check the variables
- Perform the action
- Check the action

Note: The start and stop scripts are required to be *idempotent*; that is, they have the appearance of being run once but can in fact be run multiple times. For example, if the start script is run for a resource that is already started, the script must not return an error.

All action scripts must return the status to the following file:

```
/var/cluster/ha/log/script_nodename
```

Examples of Action Scripts

The following sections use portions of the NFS scripts as examples.

Note: The examples in this guide may not exactly match the released system.

start Script

The NFS start script does the following:

1. Creates a resource-specific NFS status directory.
2. Exports the specified export-point with the specified export-options.

Following is a section from the NFS start script:

```
# Start the resource on the local machine.
# Return HA_SUCCESS if the resource has been successfully started on the local
# machine and HA_CMD_FAILED otherwise.
#
start_nfs()
{
  ${HA_DBGLOG} "Entry: start_nfs()";

  # for all nfs resources passed as parameter
  for resource in ${HA_RES_NAMES}
  do
    NFSFILEDIR=${HA_SCRIPTTMPDIR}/${LOCAL_TEST_KEY}$resource
    HA_CMD="/sbin/mkdir -p $NFSFILEDIR";
    ha_execute_cmd "creating nfs status file directory";
  done
}
```

```

if [ $? -ne 0 ]; then
    ${HA_LOG} "Failed to create ${NFSFILEDIR} directory";
    ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
    exit_script $HA_NOCFGINFO
fi

get_nfs_info $resource
if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: $resource parameters not present in CDB";
    ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
    exit_script ${HA_NOCFGINFO};
fi

ha_get_field "${HA_STRING}" export-info
if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: export-info not present in CDB for resource $resource";
    ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
    exit_script ${HA_NOCFGINFO};
fi
export_opts="$HA_FIELD_VALUE"

ha_get_field "${HA_STRING}" filesystem
if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: filesystem-info not present in CDB for resource
$resource";
    ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
    exit_script ${HA_NOCFGINFO};
fi
filesystem="$HA_FIELD_VALUE"
# Make the script idempotent, check to see if the NFS resource
# is already exported, if so return success. Remember that we
# might not have any export options.
retstat=0;
# Check to see if the NFS resource is already exported
# (without options)
/usr/etc/exportfs | grep "$resource$" >/dev/null 2>&1
retstat=$?
if [ $retstat -eq 1 ]; then
    # Check to see if the NFS resource is already exported
    # with options.
    /usr/etc/exportfs | grep "$resource " | grep "$export_opts$" >/dev/null 2>&1

```

```
    retstat=$?
fi
if [ $retstat -eq 1 ]; then
    # Before we try and export the NFS resource, make sure
    # filesystem is mounted.
    HA_CMD="/sbin/grep $filesystem /etc/mtab > /dev/null 2>&1";
    ha_execute_cmd "check if the filesystem $filesystem is mounted";
    if [ $? -eq 0 ]; then
        HA_CMD="/usr/etc/exportfs -i -o $export_opts $resource";
        ha_execute_cmd "export $resource directories to NFS clients";
        if [ $? -ne 0 ]; then
            ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
        else
            ha_write_status_for_resource ${resource} ${HA_SUCCESS};
        fi
    else
        ${HA_LOG} "NFS: filesystem $filesystem not mounted"
        ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
    fi
else
    ha_write_status_for_resource ${resource} ${HA_SUCCESS};
fi
done
}
```

stop Script

The NFS stop script does the following:

1. Unexports the specified export-point.
2. Removes the NFS status directory.

Following is an example from the NFS stop script:

```
# Stop the nfs resource on the local machine.
# Return HA_SUCCESS if the resource has been successfully stopped on the local
# machine and HA_CMD_FAILED otherwise.
#
stop_nfs()
{
```

```

${HA_DBGLOG} "Entry: stop_nfs()";

# for all nfs resources passed as parameter
for resource in ${HA_RES_NAMES}
do
get_nfs_info $resource
if [ $? -ne 0 ]; then
    # NFS resource information not available.
    ${HA_LOG} "NFS: $resource parameters not present in CDB";
    ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
    exit_script ${HA_NOCFGINFO};
fi

ha_get_field "${HA_STRING}" export-info
if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: export-info not present in CDB for resource $resource";
    ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
    exit_script ${HA_NOCFGINFO};
fi
export_opts="${HA_FIELD_VALUE}"

# Make the script idempotent, check to see if the filesystem
# is already exported, if so return success. Remember that we
# might not have any export options.

retstat=0;
# Check to see if the filesystem is already exported
# (without options)
/usr/etc/exportfs | grep "$resource$" >/dev/null 2>&1
retstat=$?
if [ $retstat -eq 1 ]; then
    # Check to see if the filesystem is already exported
    # with options.
    /usr/etc/exportfs | grep "$resource " | grep "$export_opts$" >/dev/null 2>&1
    retstat=$?
fi
if [ $retstat -eq 0 ]; then
    # Before we unexport the filesystem, check that it exists
    HA_CMD="/sbin/grep $resource /etc/mstab > /dev/null 2>&1";
    ha_execute_cmd "check if the export-point exists";
    if [ $? -eq 0 ]; then

```

```
HA_CMD="/usr/etc/exportfs -u $resource";
ha_execute_cmd "unexport $resource directories to NFS clients";
if [ $? -ne 0 ]; then
    ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
else
    ha_write_status_for_resource ${resource} ${HA_SUCCESS};
fi
else
    ${HA_LOG} "NFS: filesystem $resource not found in export filesystem list, \
unexporting anyway";
    HA_CMD="/usr/etc/exportfs -u $resource";
    ha_execute_cmd "unexport $resource directories to NFS clients";
    ha_write_status_for_resource ${resource} ${HA_SUCCESS};
fi
else
    ha_write_status_for_resource ${resource} ${HA_SUCCESS};
fi

# remove the monitor nfs status file
NFSFILEDIR=${HA_SCRIPTTMPDIR}/${LOCAL_TEST_KEY}$resource
HA_CMD="/sbin/rm -rf $NFSFILEDIR";
ha_execute_cmd "removing nfs status file directory";
if [ $? -ne 0 ]; then
    ${HA_LOG} "Failed to delete ${NFSFILEDIR} directory";
    ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
    exit_script $HA_NOCFGINFO
fi
done
}
```

monitor Script

The NFS monitor script does the following:

1. Verifies that the file system is mounted at the correct mount point.
2. Requests the status of the exported file system.
3. Checks the export-point.
4. Requests NFS statistics and (based on the results) make a Remote Procedure Call (RPC) to NFS as needed.

Following is an example from the NFS monitor script:

```
# Check if the nfs resource is allocated in the local node
# This check must be light weight and less intrusive compared to
# exclusive check. This check is done when the resource has been
# allocated in the local node.
# Return HA_SUCCESS if the resource is running in the local node
# and HA_CMD_FAILED if the resource is not running in the local node
# The list of the resources passed as input is in variable
# $HA_RES_NAMES
#
monitor_nfs()
{
  ${HA_DBGLOG} "Entry: monitor_nfs()";

  for resource in ${HA_RES_NAMES}
  do
    get_nfs_info $resource
    if [ $? -ne 0 ]; then
      # No resource information available.
      ${HA_LOG} "NFS: $resource parameters not present in CDB";
      ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
      exit_script ${HA_NOCFGINFO};
    fi

    ha_get_field "${HA_STRING}" filesystem
    if [ $? -ne 0 ]; then
      # filesystem not available available.
      ${HA_LOG} "NFS: filesystem not present in CDB for resource $resource";
      ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
      exit_script ${HA_NOCFGINFO};
    fi
  fi

  fs="${HA_FIELD_VALUE}";

  # Check to see if the filesystem is mounted
  HA_CMD="/sbin/mount | grep $fs >> /dev/null 2>&1"
  ha_execute_cmd "check to see if $fs is mounted"
  if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: $fs not mounted";
    ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
    exit_script $HA_CMD_FAILED;
  fi
}
```

2: Writing the Action Scripts and Adding Monitoring Agents

```
fi

# stat the filesystem
HA_CMD="/sbin/stat $resource";
ha_execute_cmd "stat mount point $resource"
if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: cannot stat $resource NFS export point";
    ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
    exit_script $HA_CMD_FAILED;
fi

# check the filesystem is exported
EXPORTFS="${HA_SCRIPTTMPDIR}/exportfs.$$"
/usr/etc/exportfs > $EXPORTFS 2>&1
HA_CMD="awk '{print \$1}' $EXPORTFS | grep $resource"
ha_execute_cmd " check the filesystem $resource is exported"
if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: failed to find $resource in exported filesystem list:-"
    ${HA_LOG} "`/sbin/cat ${EXPORTFS}`"
    rm -f $EXPORTFS;
    ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
    exit_script $HA_CMD_FAILED;
fi

rm -f $EXPORTFS

# create a file to hold the nfs stats. This will will be
# deleted in the stop script.
NFSFILE=${HA_SCRIPTTMPDIR}/${LOCAL_TEST_KEY}$resource/.nfsstat
NFS_STAT='/usr/etc/nfsstat -rs | /usr/bin/tail -1 | /usr/bin/awk '{print $1}'`
if [ ! -f $NFSFILE ]; then
    ${HA_LOG} "NFS: creating stat file $NFSFILE";
    echo $NFS_STAT > $NFSFILE;
    if [ $NFS_STAT -eq 0 ];then
        # do some rpcinfo's
        exec_rpcinfo;
        if [ $? -ne 0 ]; then
            ${HA_LOG} "NFS: exec_rpcinfo failed (1)";
            ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
            exit_script $HA_CMD_FAILED;
        fi
    fi
fi
```

```

    fi
else
    OLD_STAT=`/sbin/cat $NFSFILE`
    if test "X${NFS_STAT}" = "X"; then
    ${HA_LOG} "NFS: NFS_STAT is not set, reset to zero";
    NFS_STAT=0;
    fi
    if test "X${OLD_STAT}" = "X"; then
    ${HA_LOG} "NFS: OLD_STAT is not set, reset to zero";
    OLD_STAT=0;
    fi
    if [ $NFS_STAT -gt $OLD_STAT ]; then
echo $NFS_STAT > $NFSFILE;
    else
echo $NFS_STAT > $NFSFILE;
exec_rpcinfo;
if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: exec_rpcinfo failed (2)";
    ha_write_status_for_resource $resource ${HA_CMD_FAILED};
    exit_script $HA_CMD_FAILED;
fi
fi
fi
ha_write_status_for_resource $resource $HA_SUCCESS;
done
}

```

exclusive Script

The NFS exclusive script determines whether the file system is already exported. The check made by an exclusive script can be more expensive than a monitor check. FailSafe uses this script to determine if resources are running on a node in the cluster, and to thereby prevent starting resources on multiple nodes in the cluster.

Following is an example from the NFS exclusive script:

```
# Check if the nfs resource is running in the local node. This check can
# more intrusive than the monitor check. This check is used to determine
# if the resource has to be started on a machine in the cluster.
# Return HA_NOT_RUNNING if the resource is not running in the local node
# and HA_RUNNING if the resource is running in the local node
# The list of nfs resources passed as input is in variable
# $HA_RES_NAMES
#
exclusive_nfs()
{

    ${HA_DBGLOG} "Entry: exclusive_nfs()";

    # for all resources passed as parameter
    for resource in ${HA_RES_NAMES}
    do
        get_nfs_info $resource
        if [ $? -ne 0 ]; then
            # No resource information available
            ${HA_LOG} "NFS: $resource parameters not present in CDB";
            ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
            exit_script ${HA_NOCFGINFO};
        fi

        SMFILE=${HA_SCRIPTTMPDIR}/showmount.$$
        /etc/showmount -x >> ${SMFILE};
        HA_CMD="/sbin/grep $resource ${SMFILE} >> /dev/null 2>&1"
        ha_execute_cmd "checking for $resource exported directory"
        if [ $? -eq 0 ];then
            ha_write_status_for_resource ${resource} ${HA_RUNNING};
            ha_print_exclusive_status ${resource} ${HA_RUNNING};
        else
            ha_write_status_for_resource ${resource} ${HA_NOT_RUNNING};
            ha_print_exclusive_status ${resource} ${HA_NOT_RUNNING};
        fi
        rm -f ${SMFILE}
    done
}
```

restart Script

The NFS restart script exports the specified export-point with the specified export-options.

Following is an example from the restart script for NFS:

```
# Restart nfs resource
# Return HA_SUCCESS if nfs resource failed over successfully or
# return HA_CMD_FAILED if nfs resource could not be failed over locally.
# Return HA_NOT_SUPPORTED if local restart is not supported for nfs
# resource type.
# The list of nfs resources passed as input is in variable
# $HA_RES_NAMES
#
restart_nfs()
{
  ${HA_DBGLOG} "Entry: restart_nfs()";

  # for all nfs resources passed as parameter
  for resource in ${HA_RES_NAMES}
  do
    get_nfs_info $resource
    if [ $? -ne 0 ]; then
      ${HA_LOG} "NFS: $resource parameters not present in CDB";
      ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
      exit_script ${HA_NOCFGINFO};
    fi

    ha_get_field "${HA_STRING}" export-info
    if [ $? -ne 0 ]; then
      ${HA_LOG} "NFS: export-info not present in CDB for resource $resource";
      ha_write_status_for_resource ${resource} ${HA_NOCFGINFO};
      exit_script ${HA_NOCFGINFO};
    fi
    export_opts="${HA_FIELD_VALUE}"

    HA_CMD="/usr/etc/exportfs -i -o $export_opts $resource";
    ha_execute_cmd "export $resource directories to NFS clients";
    if [ $? -ne 0 ]; then
      ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
    else

```

```
    ha_write_status_for_resource ${resource} ${HA_SUCCESS};  
fi  
done  
}
```

Monitoring Agents

If resources cannot be monitored using a lightweight check, you should use a *monitoring agent*. The `monitor` action script contacts the monitoring agent to determine the status of the resource in the node. The monitoring agent in turn periodically monitors the resource. Figure 2-1 shows the monitoring process.



Figure 2-1 Monitoring Process

Monitoring agents are useful for monitoring database resources. In databases, creating the database connection is costly and time consuming. The monitoring agent maintains connections to the database and it queries the database using the connection in response to the `monitor` action script request.

Monitoring agents are independent processes and can be started by the `cmond` process, although this is not required. For example, if a monitoring agent must be started when activating highly available services on a node, information about that agent can be added to the `cmond` configuration on that node. The `cmond` configuration is located in the `/var/cluster/cmon/process_groups` directory. Information about different agents should go into different files. The name of the file is not relevant to the activate/deactivate procedure.

If a monitoring agent exits or aborts, `cmond` will automatically restart the monitoring agent. This prevents `monitor` action script failures due to monitoring agent failures.

For example, the `/var/cluster/cmon/process_groups/ip_addresses` file contains information about the `ha_ifd` process that monitors network interfaces. It contains the following:

```
TYPE = cluster_agent
PROCS = ha_ifd
ACTIONS = start stop restart attach detach
AUTOACTION = attach
```

Note: The `ACTIONS` line above defines what `cmond` can do to the `PROCS` processes. These actions must be the same for every agent. (It does not refer to action scripts.)

If you create a new monitoring agent, you must also create a corresponding file in the `/var/cluster/cmon/process_groups` directory that contains similar information about the new agent. To do this, you can copy the `ip_addresses` file and modify the `PROCS` line to list the executables that constitute your new agent. These executables must be located in the `/usr/cluster/bin` directory. You should not modify the other configuration lines (`TYPE`, `ACTIONS`, and `AUTOACTION`).

Suppose you need to add a new agent called `newagent` that consists of processes `ha_x` and `ha_y`. The configuration information for this agent will be located in the `/var/cluster/cmon/process_groups/newagent` file, which will contain the following:

```
TYPE = cluster_agent
PROCS = ha_x ha_y
ACTIONS = start stop restart attach detach
AUTOACTION = attach
```

In this case, the software will expect two executables (`/usr/cluster/bin/ha_x` and `/usr/cluster/bin/ha_y`) to be present.

Creating a Failover Policy

This chapter tells you how to create a failover policy. It describes the following topics:

- "Contents of a Failover Policy"
- "Failover Script Interface" on page 48
- "Creating a Failover Policy that Returns the Resource Group to the Same Node" on page 49
- "Example Failover Policies" on page 49

Contents of a Failover Policy

A failover policy is the method by which a resource group is failed over from one node to another. A failover policy consists of the following:

- Failover domain
- Failover attribute
- Failover scripts

FailSafe uses the failover domain output from a failover script along with failover attributes to determine on which node a resource group should reside.

The administrator must configure a failover policy for each resource group. The name of the failover policy must be unique within the pool.

Failover Domain

A failover domain is the **ordered** list of nodes on which a given resource group can be allocated. The nodes listed in the failover domain **must** be within the same cluster; however, the failover domain does not have to include every node in the cluster. The failover domain can also be used to statically load balance the resource groups in a cluster.

Examples:

- In a four-node cluster, a set of two nodes that have access to a particular XLV volume may be the failover domain of the resource group containing that XLV volume.
- In a cluster of nodes named `venus`, `mercury`, and `pluto`, you could configure the following initial failover domains for resource groups `RG1` and `RG2`:
 - `mercury, venus, pluto` for `RG1`
 - `pluto, mercury` for `RG2`

The administrator defines the initial failover domain when configuring a failover policy. The initial failover domain is used when a cluster is first booted. The ordered list specified by the initial failover domain is transformed into a run-time failover domain by the failover script. With each failure, the failover script takes the current run-time failover domain and potentially modifies it (for the `ordered` failover script, the order will not change); the initial failover domain is never used again. Depending on the run-time conditions such as load and contents of the failover script, the initial and run-time failover domains may be identical.

For example, suppose that the cluster contains three nodes named `N1`, `N2`, and `N3`; that node failure is not the reason for failover; and that the initial failover domain is as follows:

`N1 N2 N3`

The runtime failover domain will vary based on the failover script:

- If `ordered`:

`N1 N2 N3`
- If `round-robin`:

`N2 N3 N1`
- If a customized failover script, the order could be any permutation, based on the contents of the script:

<code>N1 N2 N3</code>	<code>N1 N3 N2</code>
<code>N2 N1 N3</code>	<code>N2 N3 N1</code>
<code>N3 N1 N2</code>	<code>N3 N2 N1</code>

FailSafe stores the run-time failover domain and uses it as input to the next failover script invocation.

Failover Attributes

A *failover attribute* is a value that is passed to the failover script and used by FailSafe for the purpose of modifying the run-time failover domain used for a specific resource group.

You can specify the following classes of failover attributes:

- Required attributes: either `Auto_Failback` or `Controlled_Failback` (mutually exclusive)
- Optional attributes:
 - `Auto_Recovery` or `InPlace_Recovery` (mutually exclusive)
 - `Critical_RG`
 - `Node_Failures_Only`

Note: The starting conditions for the attributes differs by class:

- For required attributes, a node joins the FailSafe membership when the cluster is already providing highly available services.
 - For optional attributes, highly available services are started and the resource group is running in only one node in the cluster.
-

Table 3-1 describes each attribute.

Table 3-1 Failover Attributes

Class	Name	Description
Required	Auto_Failback	Specifies that the resource group is made online based on the failover policy when the node joins the cluster. This attribute is best used when some type of load balancing is required. You must specify either this attribute or the <code>Controlled_Failback</code> attribute.
	Controlled_Failback	Specifies that the resource group remains on the same node when a node joins the cluster. This attribute is best used when client/server applications have expensive recovery mechanisms, such as for databases or applications that use <code>tcp</code> to communicate. You must specify either this attribute or the <code>Auto_Failback</code> attribute.
Optional	Auto_Recovery	Specifies that the resource group is made online based on the failover policy even when an exclusivity check shows that the resource group is running on a node. This attribute is optional and is mutually exclusive with the <code>InPlace_Recovery</code> attribute. If you specify neither of these attributes, <code>FailSafe</code> will use this attribute by default if you have specified the <code>Auto_Failback</code> attribute.
	InPlace_Recovery	Specifies that the resource group is made online on the same node where the resource group is running. This attribute is optional and is mutually exclusive with the <code>Auto_Recovery</code> attribute. If you specify neither of these attributes, <code>FailSafe</code> will use this attribute by default if you have specified the <code>Controlled_Failback</code> attribute.

Class	Name	Description
	Critical_RG	Allows monitor failure recovery to succeed even when there are resource group release failures. When resource monitoring fails, FailSafe attempts to move the resource group to another node in the application failover domain. If FailSafe fails to release the resources in the resource group, FailSafe puts the Resource group into <code>srmd executable error</code> status. If the <code>Critical_RG</code> failover attribute is specified in the failover policy of the resource group, FailSafe will reset the node where the release operation failed and move the resource group to another node based on the failover policy.
	Node_Failures_Only	Allows failover only when there are node failures. This attribute does not have an impact on resource restarts in the local node. The failover does not occur when there is a resource monitoring failure in the resource group. This attribute is useful for customers who are using a hierarchical storage management system such as DMF; in this situation, a customer may want to have resource monitoring failures reported without automatic recovery, allowing operators to perform the recovery action manually if necessary.

Failover Scripts

A failover script generates the run-time failover domain and returns it to the FailSafe process. The FailSafe process applies the failover attributes and then selects the first node in the returned failover domain that is also in the current FailSafe membership.

Note: The run-time of the failover script must be capped to a system-definable maximum. Therefore, any external calls must be guaranteed to return quickly. If the failover script takes too long to return, FailSafe will kill the script process and use the previous run-time failover domain.

Failover scripts are stored in the `/var/cluster/ha/policies` directory.

ordered

The `ordered` failover script is provided with the release. The `ordered` script never changes the initial domain; when using this script, the initial and run-time domains

are equivalent. The script reads six lines from the input file and in case of errors logs the input parameters and/or the error to the script log.

The following example shows the contents of the ordered failover script:

```
#!/sbin/ksh
#
# $1 - input file
# $2 - output file
#
# line 1 input file - version
# line 2 input file - name
# line 3 input file - owner field
# line 4 input file - attributes
# line 5 input file - list of possible owners
# line 6 input file - application failover domain

DIR=/usr/cluster/bin
LOG=${DIR}/ha_cilog -g ha_script -s script
FILE=/var/cluster/ha/policies/ordered

input=$1
output=$2
cat ${input} | read version
head -2 ${input} | tail -1 | read name
head -3 ${input} | tail -1 | read owner
head -4 ${input} | tail -1 | read attr
head -5 ${input} | tail -1 | read mem1 mem2 mem3 mem4 mem5 mem6 mem7 mem8
head -6 ${input} | tail -1 | read afd1 afd2 afd3 afd4 afd5 afd6 afd7 afd8

${LOG} -l 1 "${FILE}:" ` /bin/cat ${input} `

if [ "${version}" -ne 1 ] ; then
    ${LOG} -l 1 "ERROR: ${FILE}: Different version no. Should be (1) rather than
    (${version})" ;
    exit 1;
elif [ -z "${name}" ] ; then
    ${LOG} -l 1 "ERROR: ${FILE}: Failover script not defined";
    exit 1;
elif [ -z "${attr}" ] ; then
    ${LOG} -l 1 "ERROR: ${FILE}: Attributes not defined";
    exit 1;
```

```

elif [ -z "${mem1}" ]; then
    ${LOG} -l 1 "ERROR: ${FILE}: No node membership defined";
    exit 1;
elif [ -z "${afd1}" ]; then
    ${LOG} -l 1 "ERROR: ${FILE}: No failover domain defined";
    exit 1;
fi

found=0
for i in $afd1 $afd2 $afd3 $afd4 $afd5 $afd6 $afd7 $afd8; do
    for j in $mem1 $mem2 $mem3 $mem4 $mem5 $mem6 $mem7 $mem8; do
        if [ "X${j}" = "X${i}" ]; then
            found=1;
            break;
        fi
    done
done

if [ ${found} -eq 0 ]; then
    mem=("${mem1}" "${mem2}" "${mem3}" "${mem4}" "${mem5}"
    "${mem6}" "${mem7}" "${mem8}");
    afd=("${afd1}" "${afd2}" "${afd3}" "${afd4}" "${afd5}"
    "${afd6}" "${afd7}" "${afd8}");
    ${LOG} -l 1 "ERROR: ${FILE}: Policy script failed"
    ${LOG} -l 1 "ERROR: ${FILE}: " `bin/cat ${input}`
    ${LOG} -l 1 "ERROR: ${FILE}: Nodes defined in membership do not match the
ones in failure domain"
    ${LOG} -l 1 "ERROR: ${FILE}: Parameters read from input file: version =
$version, name = $name, owner = $owner, attribute = $attr, nodes = $mem, afd = $afd"
    exit 1;
fi

if [ ${found} -eq 1 ]; then
    rm -f ${output}
    echo $afd1 $afd2 $afd3 $afd4 $afd5 $afd6 $afd7 $afd8 > ${output}
    exit 0
fi
exit 1

```

round-robin

The round-robin failover script selects the resource group owner in a round-robin (circular) fashion. This policy can be used for resource groups that can be run in any node in the cluster.

The following example shows the contents of the round-robin failover script:

```
#!/sbin/ksh
#
# $1 - input file
# $2 - output file
#
# line 1 input file - version
# line 2 input file - name
# line 3 input file - owner field
# line 4 input file - attributes
# line 5 input file - Possible list of owners
# line 6 input file - application failover domain

DIR=/usr/cluster/bin
LOG=${DIR}/ha_cilog -g ha_script -s script
FILE=/var/cluster/ha/policies/round-robin

# Read input file
input=$1
output=$2
cat ${input} | read version
head -2 ${input} | tail -1 | read name
head -3 ${input} | tail -1 | read owner
head -4 ${input} | tail -1 | read attr
head -5 ${input} | tail -1 | read mem1 mem2 mem3 mem4 mem5 mem6 mem7 mem8
head -6 ${input} | tail -1 | read afd1 afd2 afd3 afd4 afd5 afd6 afd7 afd8

# Validate input file
${LOG} -1 1 "${FILE}:" ` /bin/cat ${input} `

if [ "${version}" -ne 1 ] ; then
    ${LOG} -1 1 "ERROR: ${FILE}: Different version no. Should be (1) rather than
    (${version})" ;
    exit 1;
```

```
elif [ -z "${name}" ]; then
    ${LOG} -l 1 "ERROR: ${FILE}: Failover script not defined";
    exit 1;
elif [ -z "${attr}" ]; then
    ${LOG} -l 1 "ERROR: ${FILE}: Attributes not defined";
    exit 1;
elif [ -z "${mem1}" ]; then
    ${LOG} -l 1 "ERROR: ${FILE}: No node membership defined";
    exit 1;
elif [ -z "${afd1}" ]; then
    ${LOG} -l 1 "ERROR: ${FILE}: No failover domain defined";
    exit 1;
fi

# Return 0 if $1 is in the membership and return 1 otherwise.
check_in_mem()
{
    for j in $mem1 $mem2 $mem3 $mem4 $mem5 $mem6 $mem7 $mem8; do
        if [ "X${j}" = "X$1" ]; then
            return 0;
        fi
    done
    return 1;
}

# Check if owner has to be changed. There is no need to change owner if
# owner node is in the possible list of owners.
check_in_mem ${owner}
if [ $? -eq 0 ]; then
    nextowner=${owner};
fi

# Search for the next owner
if [ "X${nextowner}" = "X" ]; then
    next=0;
    for i in $afd1 $afd2 $afd3 $afd4 $afd5 $afd6 $afd7 $afd8; do
        if [ "X${i}" = "X${owner}" ]; then
            next=1;
            continue;
        fi
    fi
```

3: Creating a Failover Policy

```
    if [ "X${owner}" = "XNO ONE" ]; then
        next=1;
    fi

    if [ ${next} -eq 1 ]; then
        # Check if ${i} is in membership
        check_in_mem ${i};
        if [ $? -eq 0 ]; then
            # found next owner
            nextowner=${i};
            next=0;
            break;
        fi
    fi
done
fi

if [ "X${nextowner}" = "X" ]; then
    # wrap round the afd list.
    for i in $afd1 $afd2 $afd3 $afd4 $afd5 $afd6 $afd7 $afd8; do
        if [ "X${i}" = "X${owner}" ]; then
            # Search for next owner complete
            break;
        fi
    fi

    # Previous loop should have found new owner
    if [ "X${owner}" = "XNO ONE" ]; then
        break;
    fi

    if [ ${next} -eq 1 ]; then
        check_in_mem ${i};
        if [ $? -eq 0 ]; then
            # found next owner
            nextowner=${i};
            next=0;
            break;
        fi
    fi
done
```

```
fi

if [ "X${nextowner}" = "X" ]; then
    ${LOG} -l 1 "ERROR: ${FILE}: Policy script failed"
    ${LOG} -l 1 "ERROR: ${FILE}: " `'/bin/cat ${input}`
    ${LOG} -l 1 "ERROR: ${FILE}: Could not find new owner"
    exit 1;
fi

# nextowner is the new owner
print=0;
rm -f ${output};

# Print the new afd to the output file
echo -n "${nextowner} " > ${output};
for i in $afd1 $afd2 $afd3 $afd4 $afd5 $afd6 $afd7 $afd8;
do
    if [ "X${nextowner}" = "X${i}" ]; then
        print=1;
    elif [ ${print} -eq 1 ]; then
        echo -n "${i} " >> ${output}
    fi
done

print=1;
for i in $afd1 $afd2 $afd3 $afd4 $afd5 $afd6 $afd7 $afd8; do
    if [ "X${nextowner}" = "X${i}" ]; then
        print=0;
    elif [ ${print} -eq 1 ]; then
        echo -n "${i} " >> ${output}
    fi
done

echo >> ${output};
exit 0;
```

Creating a New Failover Script

If the ordered or round-robin scripts do not meet your needs, you can create a new failover script and place it in the `/var/clusters/ha/policies` directory. You can then configure the cluster database to use your new failover script for the required resource groups.

Failover Script Interface

The following is passed to the failover script:

```
function(version, name, owner, attributes, possibleowners, domain)
```

<i>version</i>	FailSafe version. The IRIX FailSafe 2.1.x release uses version number 1.
<i>name</i>	Name of the failover script (used for error validations and logging purposes).
<i>owner</i>	Logical name of the node that has (or had) the resource group online.
<i>attributes</i>	Failover attributes (<code>Auto_Failback</code> or <code>Controlled_Failback</code> must be included).
<i>possibleowners</i>	List of possible owners for the resource group. This list can be a subset of the current FailSafe membership.
<i>domain</i>	Ordered list of nodes used at the last failover. (At the first failover, the initial failover domain is used.)

The failover script returns the newly generated run-time failover domain to FailSafe, which then chooses the node on which the resource group should be allocated by applying the failover attributes and FailSafe membership to the run-time failover domain.

Creating a Failover Policy that Returns the Resource Group to the Same Node

When HA services are stopped, a resource group will continue to run on the same node (that is, the node where it was running at the time HA services were stopped), as long as the node is available as part of the cluster. The resource group will switch to a different node in the failover domain only when the node is not available.

When resources are started, the node information can be stored in the configuration database as the resource group owner. When resources are stopped, the resource group owner information is removed. When FailSafe is started, the resource group ownership is read from the configuration database as part of the failover policy. This information is used by the failover policy script to determine the node where resource group should run.

Example Failover Policies

There are two general types of configuration, each of which can have from two through eight nodes:

- N nodes that can potentially failover their applications to any of the other nodes in the cluster.
- N primary nodes that can failover to M backup nodes. For example, you could have three primary nodes and one backup node.

This section shows examples of failover policies for the following types of configuration, each of which can have from two through eight nodes:

- N primary nodes and one backup node ($N+1$)
- N primary nodes and two backup nodes ($N+2$)
- N primary nodes and M backup nodes ($N+M$)

Note: The diagrams in the following sections illustrate the configuration concepts discussed here, but they do not address all required or supported elements, such as reset hubs.

N+1 Configuration

Figure 3-1 shows a specific instance of an N+1 configuration in which there are three primary nodes and one backup node. (This is also known as a *star configuration*.) The disks shown could each be disk farms.

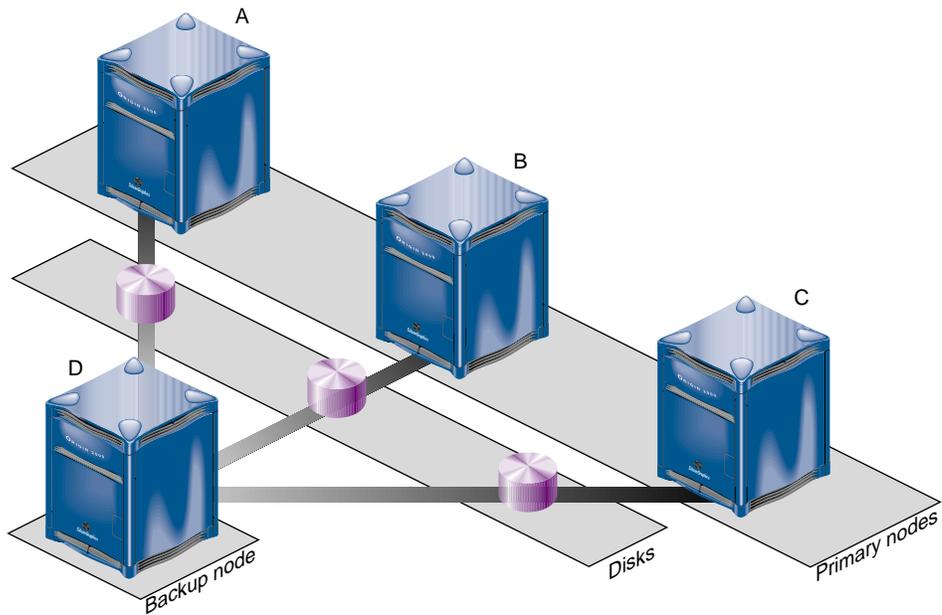


Figure 3-1 N+1 Configuration Concept

You could configure the following failover policies for load balancing:

- Failover policy for RG1:
 - Initial failover domain = A, D
 - Failover attribute = Auto_Failback, Critical_RG
 - Failover script = ordered
- Failover policy for RG2:
 - Initial failover domain = B, D

- Failover attribute = `Auto_Failback`
- Failover script = `ordered`
- Failover policy for RG3:
 - Initial failover domain = C, D
 - Failover attribute = `Auto_Failback`
 - Failover script = `ordered`

If node A fails, RG1 will fail over to node D. As soon as node A reboots, RG1 will be moved back to node A.

If you change the failover attribute to `Controlled_Failback` for RG1 and node A fails, RG1 will fail over to node D and will remain running on node D even if node A reboots.

Suppose resource group RG1 is online on node A in the cluster. When the monitor of one of the resources in RG1 fails, FailSafe attempts to move the resource group to node D. If the release of RG1 from node A fails, FailSafe will reset node A and allocate the resource group on node D. If `Critical_RG` failover attribute was not specified, RG1 will have an `srmd` executable error.

N+2 Configuration

Figure 3-2 shows a specific instance of an *N+2* configuration in which there are four primary nodes and two backup nodes. The disks shown could each be disk farms.

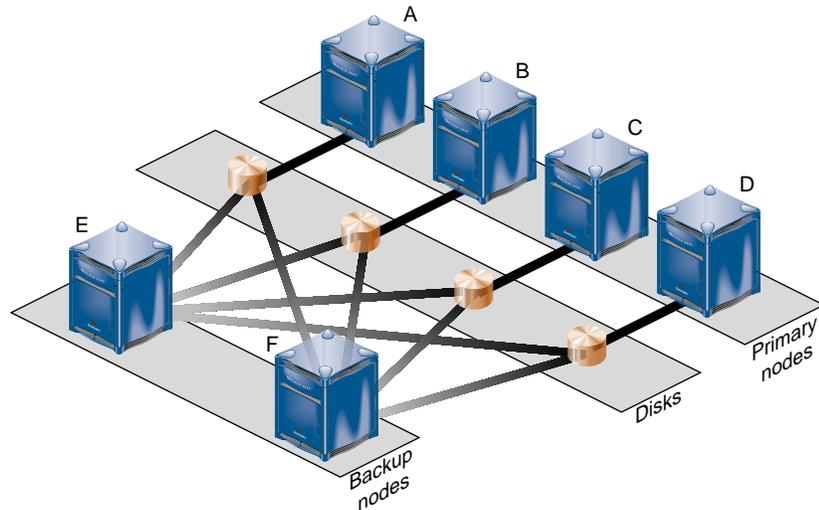


Figure 3-2 N+2 Configuration Concept

You could configure the following failover policy for resource groups RG7 and RG8:

- Failover policy for RG7:
 - Initial failover domain = A, E, F
 - Failover attribute = `Controlled_Failback`
 - Failover script = `ordered`
- Failover policy for RG8:
 - Initial failover domain = B, F, E
 - Failover attribute = `Auto_Failback`
 - Failover script = `ordered`

If node A fails, RG7 will fail over to node E. If node E also fails, RG7 will fail over to node F. If A is rebooted, RG7 will remain on node F.

If node B fails, RG8 will fail over to node F. If B is rebooted, RG8 will return to node B.

N+M Configuration

Figure 3-3 shows a specific instance of an *N+M* configuration in which there are four primary nodes and each can serve as a backup node. The disk shown could be a disk farm.

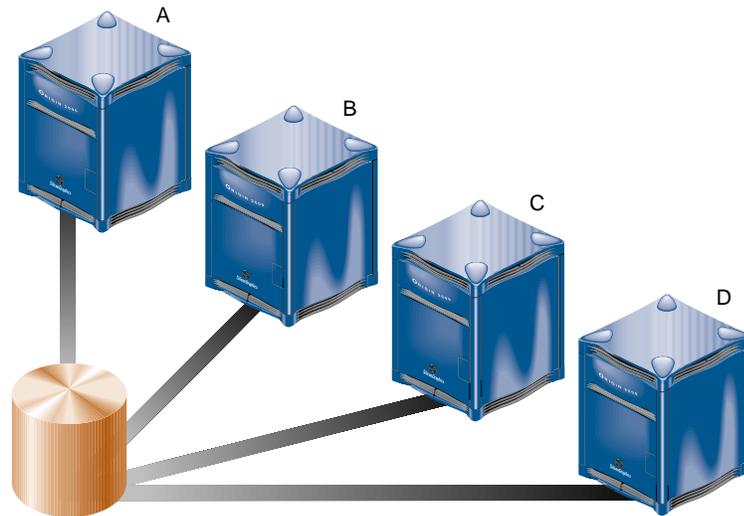


Figure 3-3 *N+M* Configuration Concept

You could configure the following failover policy for resource groups RG5 and RG6:

- Failover policy for RG5:
 - Initial failover domain = A, B, C, D
 - Failover attribute = Controlled_Failback
 - Failover script = ordered
- Failover policy for RG6:
 - Initial failover domain = C, A, D
 - Failover attribute = Controlled_Failback
 - Failover script = ordered

If node C fails, RG6 will fail over to node A. When node C reboots, RG6 will remain running on node A. If node A then fails, RG6 will return to node C and RG5 will move to node B. If node B then fails, RG5 moves to node C.

Defining a New Resource Type

Already installed resource types are created automatically when a FailSafe cluster is created. This chapter describes how to define a **new** resource type, which can be performed using the FailSafe GUI or the `cmgr` command.

The following are examples of candidates for resource types:

- Databases that support transactions
- Web servers
- Applications that require user datagram protocol (UDP) for communication with clients

See also "Characteristics that Permit an Application to be Highly Available" on page 3.

You will want to define a new resource type when creating something entirely new or when you want to have multiple resource types that are similar except for one or two attributes. For example, if you want to enable local restart for most IP addresses but not for some, you could create a new resource type called `IP_address2` using all of the same information as for the default `IP_address` except changing the value of the restart mode to 1 rather than the default 0.

This chapter contains the following sections:

- "Information You Must Gather" on page 56
- "Copying an Existing Resource Type to Create a New One" on page 59
- "Creating a New Resource Type from Scratch" on page 60
- "Server-side Properties File" on page 73
- "Testing a New Resource Type" on page 76

Information You Must Gather

To define a new resource type, you must have the following information:

- Name of the resource type. The name can consist of alphanumeric characters and any of the following:

- (hyphen) / = : @
_ (underscore) . " ,

The name cannot contain a space, an unprintable character, or any of the following characters:

* ? \ #

- Name of the cluster to which the resource type will apply.
- If the resource type is to be restricted to a specific node, you must know the node name.
- Order of performing the action scripts for resources of this type in relation to resources of other types. If multiple resource groups are being made online, the `start` action for the resource groups are done in parallel. Resources in each group are started in increasing order of their resource type order. If there are multiple resources of the same type in the resource group, they are passed as parameters to same invocation of resource type start script. For example, all resources of resource type order 10 will be started after resources with resource type order 5 in the resource group.

Ensure that the number you choose for a new resource type permits the resource types on which it depends to be started before it is started, or stopped after it is stopped, as appropriate.

Table 4-1 shows the conventions used for order ranges. The values available for customer use are 201-400 and 701-999.

Table 4-1 Order Ranges

Range	Reservation
2-100	SGI-provided basic system resource types, such as <code>MAC_address</code>
101-200	SGI-provided system plug-ins that can be started before <code>IP_address</code>
201-400	User-defined resource types that can be started before <code>IP_address</code>
401-500	SGI-provided basic system resource types, such as <code>IP_address</code>
501-700	SGI-provided system plug-ins that must be started after <code>IP_address</code>
701-999	User-defined resource types that must be started after <code>IP_address</code>

Table 4-2 shows the order numbers of the resource types provided with the release.

Table 4-2 Resource Type Order Numbers

Order Number	Resource Type
10	<code>MAC_address</code>
20	<code>volume</code>
30	<code>filesystem</code>
201	<code>NFS</code>
401	<code>IP_address</code>
412	<code>statd_unlimited</code>
500	<code>CXFS</code>
501	<code>Netscape_web</code>
502	<code>Samba</code>
511	<code>Oracle_DB</code>
521	<code>INFORMIX_DB</code>

- Restart mode, which can be one of the following values:
 - 0 = Do not restart upon monitoring failures
 - 1 = Restart a fixed number of times

- Number of local restarts (when restart mode is 1).
- Location of the executable script. This is always as follows: `/var/cluster/ha/resource_types/resource_type_tname`
- Monitoring interval, which is the time period (in milliseconds) between successive executions of the `monitor` action script; this is only valid for the `monitor` action script.
- Starting time for monitoring. When the resource group is made online in a cluster node, FailSafe will start monitoring the resources after the specified time period (in milliseconds).
- Action scripts to be defined for this resource type. You must specify scripts for `start`, `stop`, `exclusive`, and `monitor`, although the `monitor` script may contain only a return-success function if you wish. If you specify 1 for the restart mode, you must specify a `restart` script.
- Type-specific attributes to be defined for this resource type. The action scripts use this information to start, stop, and monitor a resource of this resource type. For example, NFS requires the following resource keys:
 - `export-point`, which takes a value that defines the export disk name. This name is used as input to the `exportfs` command. For example:

```
export-point = /this_disk
```
 - `export-info`, which takes a value that defines the export options for the file system. These options are used in the `exportfs` command. For example:

```
export-info = rw,wsync,anon=root
```
 - `filesystem`, which takes a value that defines the raw file system. This name is used as input to the `mount` command. For example:

```
filesystem = /dev/xlv/xlv_object
```

Copying an Existing Resource Type to Create a New One



Caution: Any user can use the GUI to **view** database information; therefore, you should not include any sensitive information in the cluster database. Users should keep this mind when deciding the list of resource attributes for the resource type.

If an existing resource type is similar to the type you want to create, you can use the following procedure:

1. Log in as `root`.
2. Copy the directory for the existing resource type and give the new directory an appropriate name. For example, to use the `NFS` resource type as the basis for a new resource type named `NFS_CXFS`, do the following:

```
# cp -r /var/cluster/ha/resource_types/NFS /var/cluster/ha/resource_types/NFS_CXFS
```

3. Modify each script in the new `NFS_CXFS` directory so that it uses the name of the new resource type. You must make this modification for the `LOCAL_TEST_KEY=` variable definition; modifying log messages and comments is optional but recommended.

For example, you would change the `LOCAL_TEST_KEY=NFS` line in the `/var/cluster/ha/resource_types/NFS_CXFS/start` script as follows:

- From:

```
LOCAL_TEST_KEY=NFS
```

- To:

```
LOCAL_TEST_KEY=NFS_CXFS
```

4. Eliminate any unneeded dependencies for the new resource type, using either the GUI or the `cmgr` command.

For example, you would eliminate the `filesystem` dependency from the new `NFS_CXFS` as follows:

```
# cmgr
Welcome to SGI Cluster Manager Command-Line Interface
cmgr> modify resource_type NFS_CXFS in cluster "testcluster"
Enter commands, when finished enter either "done" or "cancel"
resource_type NFS_CXFS ? remove dependency filesystem
```

```
resource_type NFS_CXFS ? done
Successfully modified resource_type NFS_CXFS
```

5. Modify the monitor script for the new resource type as needed.

For example, the difference between the standard NFS monitor script and the new NFS_CXFS monitor script is that when you export CXFS filesystems, you do not want FailSafe to check if the filesystem is mounted and to exit with HA_CMD_FAILED if it is not. The NFS_CXFS monitor script itself will determine what action should take place if the filesystem becomes unmounted. To accomplish this, you would modify the `/var/cluster/ha/resource_types/NFS_CXFS/monitor` script to comment out the `exit_script` line in the following section (the line as modified is shown here in bold):

```
# Check to see if the filesystem is mounted
HA_CMD="/sbin/mount | grep $fs >> /dev/null 2>&1"
ha_execute_cmd "check to see if $fs is mounted"
if [ $? -ne 0 ]; then
    ${HA_LOG} "NFS: $fs not mounted";
    ha_write_status_for_resource ${resource} ${HA_CMD_FAILED};
# exit_script $HA_CMD_FAILED;
fi
```

The result of this change is that the status of the commands will be written to the log, but the script will not exit.

Creating a New Resource Type from Scratch



Caution: Any user can use the GUI to **view** database information; therefore, you should not include any sensitive information in the cluster database. Users should keep this mind when deciding the list of resource attributes for the resource type.

If none of the existing resource types are similar to the type you want to create, you can create a resource type from scratch using the following methods:

- "Using the FailSafe Manager GUI " on page 61
- "Using `cmgr` Interactively " on page 66
- "Using `cmgr` With a Script" on page 71

Using the FailSafe Manager GUI

You can use the FailSafe Manager graphical user interface (GUI) to define a new resource type and to define the dependencies for a given type. For details about the GUI, see Appendix A, "Starting the FailSafe Manager GUI" on page 87 and the *FailSafe Administrator's Guide for SGI InfiniteStorage*.

Define a New Resource Type

To define a new resource type using the GUI, select the following menu:

- Tasks**
 - > Resource Types**
 - > Define a Resource Type**

The GUI will prompt you for required and optional information. Online help is provided for each item.

The following figures show this process for a new resource type called `newresourcetype`.

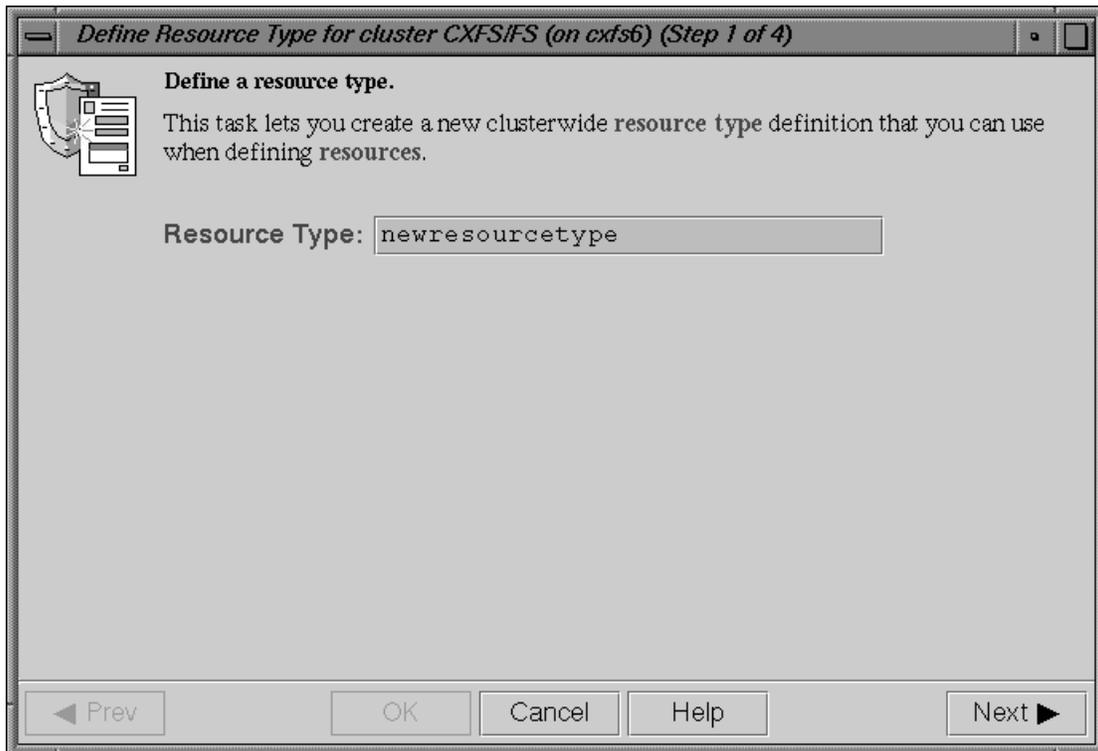


Figure 4-1 Specify the Name of the New Resource Type

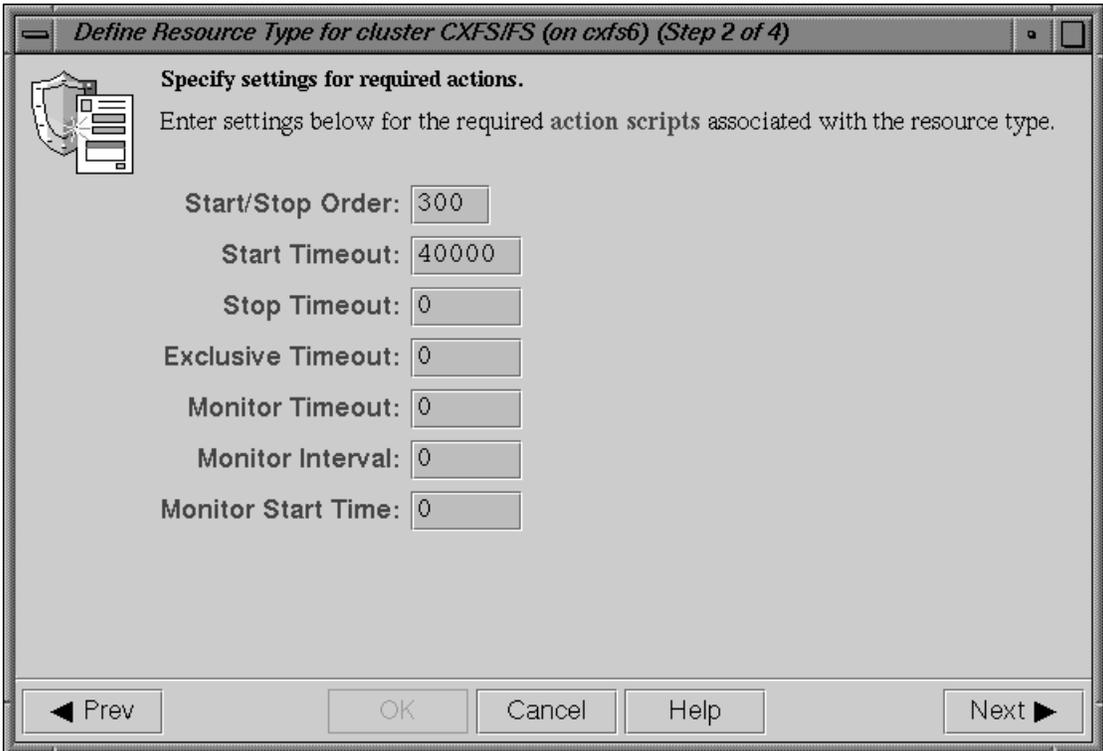


Figure 4-2 Specify Settings for Required Actions

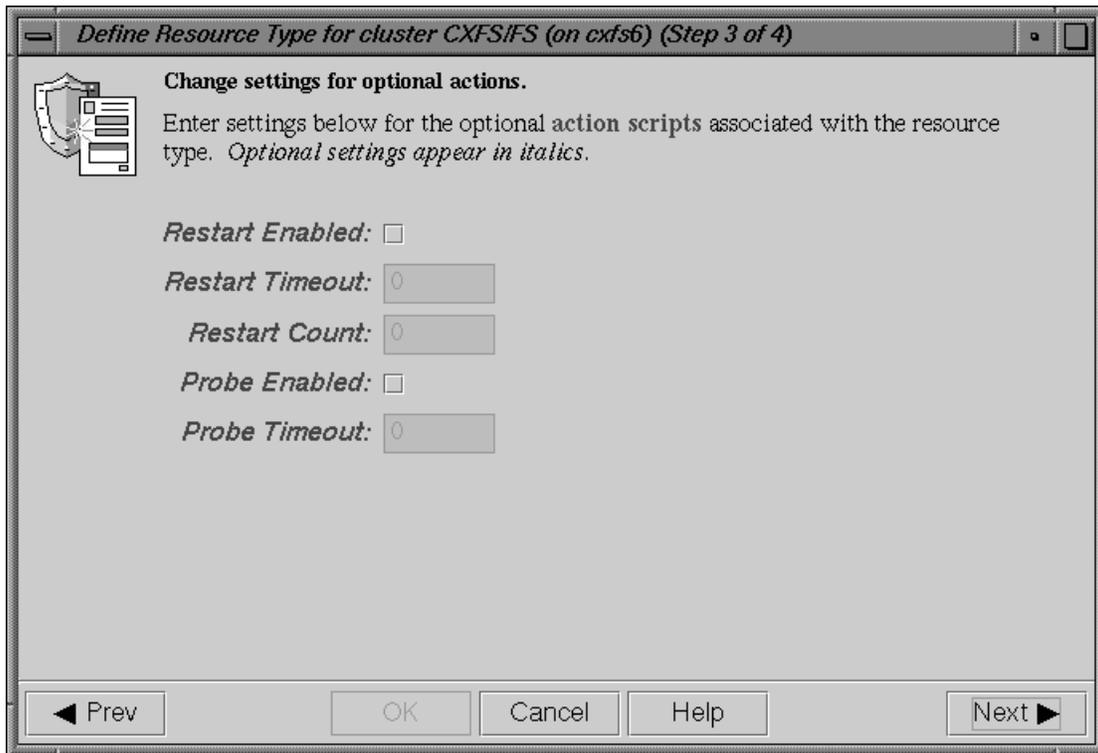


Figure 4-3 Change Settings for Optional Actions

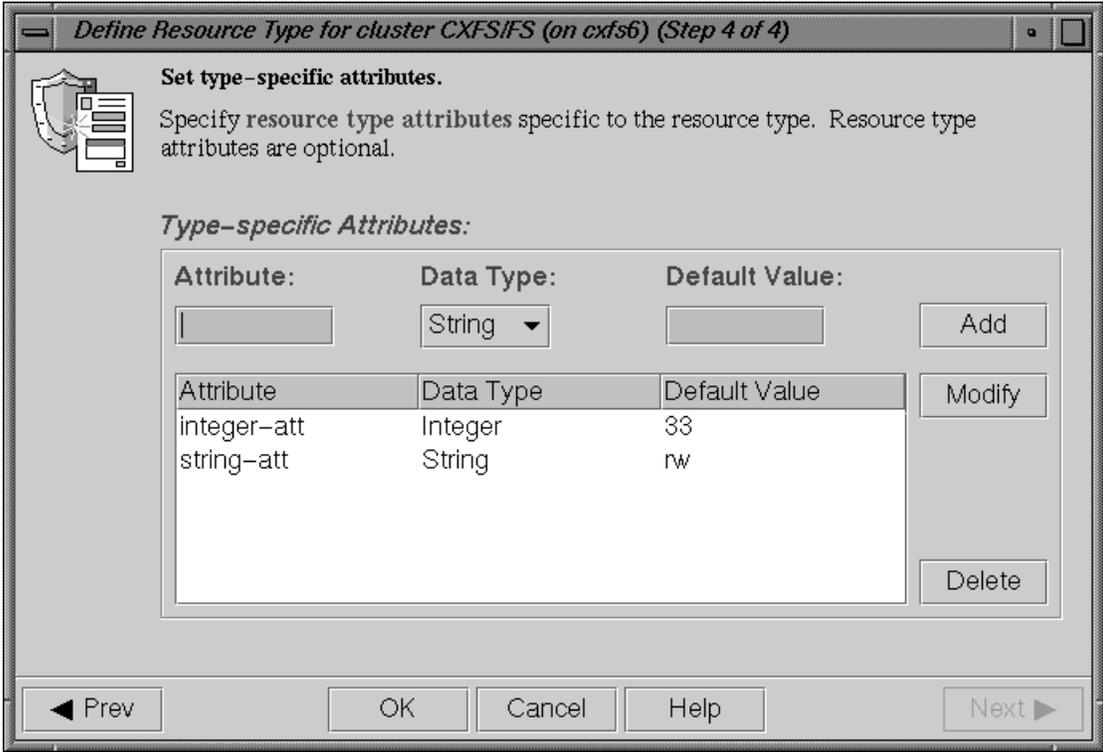


Figure 4-4 Set Type-specific Attributes

Define Dependencies

To define the dependencies for a given type, select the following menu:

- Tasks
 - > Resource Types
 - > Add/Remove Dependencies for a Resource Type

Figure 4-5 shows an example of adding a dependency (filesystem) to the newresourcetype resource type.

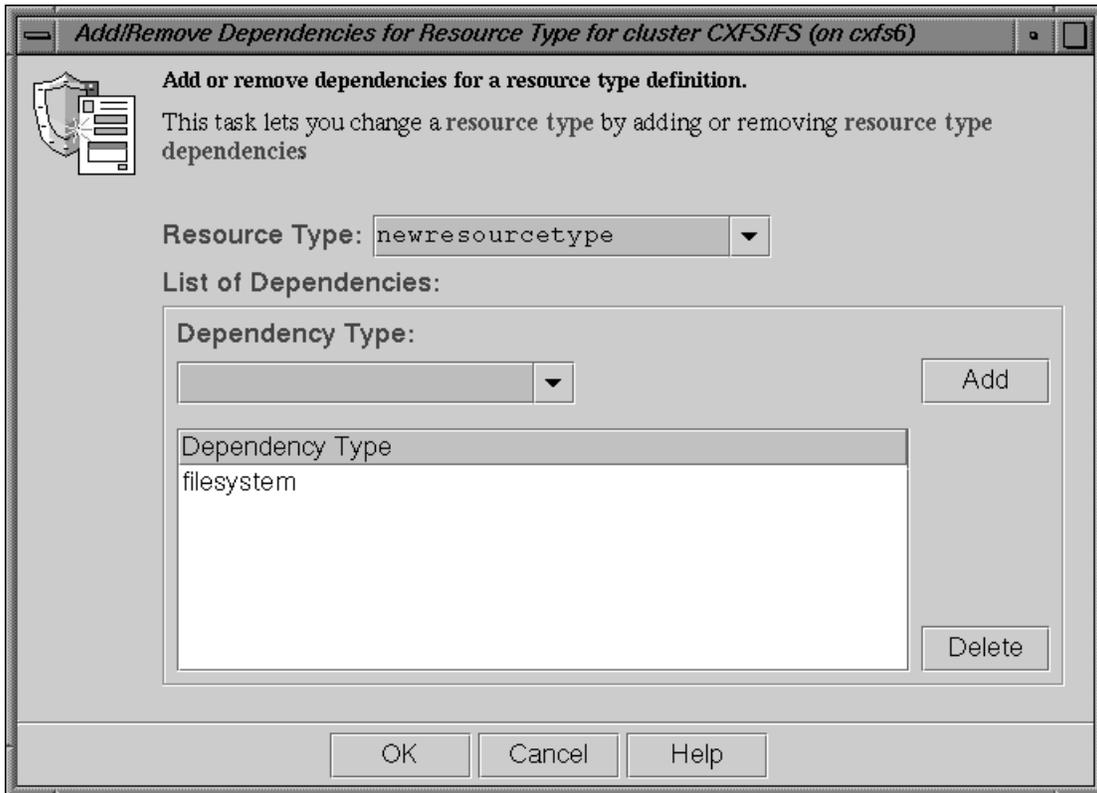


Figure 4-5 Add Dependencies

Using cmgr Interactively

The following steps show the use of `cmgr` interactively to define a resource type called `newresourcetype`.

Note: A resource type name cannot contain a space, an unprintable character, or any of the following characters:

* ? \ #

1. Log in as root.

2. Execute the `cmgr` command. You can use the `-p` option to prompt you for information:

```
# /usr/cluster/bin/cmgr -p
Welcome to SGI Cluster Manager Command-Line Interface

cmgr>
```

3. Use the `set` subcommand to specify the default cluster used for `cmgr` operations. In this example, we use a cluster named `TEST`:

```
cmgr> set cluster TEST
```

If you prefer, you can specify the cluster name as needed with each subcommand.

4. Use the `define resource_type` subcommand. By default, the resource type will apply across the cluster; if you wish to limit the resource type to a specific node, enter the node name when prompted. If you wish to enable restart mode, enter 1 when prompted.

Note: The following example only shows the prompts and answers for two action scripts (`start` and `stop`) for a new resource type named `newresourcetype`.

```
cmgr> define resource_type newresourcetype
```

```
(Enter "cancel" at any time to abort)
```

```
Node[optional]?
```

```
Order ? 300
```

```
Restart Mode ? (0)
```

```
DEFINE RESOURCE TYPE OPTIONS
```

- 0) Modify Action Script.
- 1) Add Action Script.
- 2) Remove Action Script.
- 3) Add Type Specific Attribute.
- 4) Remove Type Specific Attribute.
- 5) Add Dependency.
- 6) Remove Dependency.
- 7) Show Current Information.

4: Defining a New Resource Type

- 8) Cancel. (Aborts command)
- 9) Done. (Exits and runs command)

Enter option:1

No current resource type actions

Action name ? **start**

Executable timeout (in milliseconds) ? **40000**

- 0) Modify Action Script.
- 1) Add Action Script.
- 2) Remove Action Script.
- 3) Add Type Specific Attribute.
- 4) Remove Type Specific Attribute.
- 5) Add Dependency.
- 6) Remove Dependency.
- 7) Show Current Information.
- 8) Cancel. (Aborts command)
- 9) Done. (Exits and runs command)

Enter option:1

Current resource type actions:

start

Action name **stop**

Executable timeout? (in milliseconds) **40000**

- 0) Modify Action Script.
- 1) Add Action Script.
- 2) Remove Action Script.
- 3) Add Type Specific Attribute.
- 4) Remove Type Specific Attribute.
- 5) Add Dependency.
- 6) Remove Dependency.
- 7) Show Current Information.
- 8) Cancel. (Aborts command)
- 9) Done. (Exits and runs command)

Enter option:3

No current type specific attributes

Type Specific Attribute ? **integer-att**

Datatype ? **integer**

Default value[optional] ? **33**

- 0) Modify Action Script.
- 1) Add Action Script.
- 2) Remove Action Script.
- 3) Add Type Specific Attribute.
- 4) Remove Type Specific Attribute.
- 5) Add Dependency.
- 6) Remove Dependency.
- 7) Show Current Information.
- 8) Cancel. (Aborts command)
- 9) Done. (Exits and runs command)

Enter option:**3**

Current type specific attributes:

Type Specific Attribute - 1: integer-att

Type Specific Attribute ? **string-att**

Datatype ? **string**

Default value[optional] ? **rw**

- 0) Modify Action Script.
- 1) Add Action Script.
- 2) Remove Action Script.
- 3) Add Type Specific Attribute.
- 4) Remove Type Specific Attribute.
- 5) Add Dependency.
- 6) Remove Dependency.
- 7) Show Current Information.
- 8) Cancel. (Aborts command)
- 9) Done. (Exits and runs command)

Enter option:**5**

No current resource type dependencies

4: Defining a New Resource Type

Dependency name ? **filesystem**

- 0) Modify Action Script.
- 1) Add Action Script.
- 2) Remove Action Script.
- 3) Add Type Specific Attribute.
- 4) Remove Type Specific Attribute.
- 5) Add Dependency.
- 6) Remove Dependency.
- 7) Show Current Information.
- 8) Cancel. (Aborts command)
- 9) Done. (Exits and runs command)

Enter option:7

Current resource type actions:

- Action - 1: start
- Action - 2: stop

Current type specific attributes:

- Type Specific Attribute - 1: integer-att
- Type Specific Attribute - 2: string-att

No current resource type dependencies

Resource dependencies to be added:

- Resource dependency - 1: filesystem

- 0) Modify Action Script.
- 1) Add Action Script.
- 2) Remove Action Script.
- 3) Add Type Specific Attribute.
- 4) Remove Type Specific Attribute.
- 5) Add Dependency.
- 6) Remove Dependency.
- 7) Show Current Information.
- 8) Cancel. (Aborts command)
- 9) Done. (Exits and runs command)

Enter option:9

```
Successfully defined resource_type newresourcetype
```

```
cmgr> show resource_types
```

```
template  
MAC_address  
newresourcetype  
IP_address  
filesystem  
volume
```

```
cmgr> exit
```

```
#
```

Using cmgr With a Script

You can write a script that contains all of the information required to define a resource type and supply it to cmgr by using the `-f` option:

```
cmgr -f scriptname
```

Or, you could include the following as the first line of the script and then execute the script itself:

```
#!/usr/cluster/bin/cmgr -f
```

If any line of the script fails, `cmgr` will exit. You can choose to ignore the failure and continue the process by using the `-i` option, as follows:

```
#!/usr/cluster/bin/cmgr -if
```

Note: If you include `-i` when using a `cmgr` command line as the first line of the script, you must use this exact syntax (that is, `-if`).

A template script for creating a new resource type is located in `/var/cluster/cmgr-templates/cmgr-create-resource_type`. Each line of the script must be a valid `cmgr` line, a comment line (starting with `#`), or a blank line. You must include a `done` command line to finish a multilevel command. If you concatenate information from multiple template scripts to prepare your cluster configuration, you must remove the `quit` at the end of each template script.

For example, you could use the following script to define the same `newresourcetype` resource type defined interactively in the previous section:

```
# Script to define the "newresourcetype" resource type

set cluster TEST
define resource_type newresourcetype
set order to 300
set restart_mode to 0

add action start
set exec_time to 40000
done

add action stop
set exec_time to 40000
done

add type_attribute integer-att
set data_type to integer
set default_value to 33
done

add type_attribute string-att
```

```
set data_type to string
set default_value to rw
done
```

```
add dependency filesystem
done
```

```
quit
```

When you execute the `cmgr -f` command line with this script, you will see the following output:

```
# /usr/cluster/bin/cmgr -f newresourcetype.script
Successfully defined resource_type newresourcetype
```

To verify that the resource type was defined, enter the following:

```
# /usr/cluster/bin/cmgr -c "show resource_types in cluster TEST"
```

```
template
MAC_address
newresourcetype
IP_address
filesystem
volume
```

Server-side Properties File

Each resource type can have an optional properties file containing a formatted label for each plug-in attribute and strings of help text that will be displayed in the GUI. The file has the following name:

```
/var/cluster/ha/resource_types/resource_type/resource_type
```

For example, the properties file for the `IP_Address` resource type would be as follows:

```
/var/cluster/ha/resource_types/IP_Address/IP_Address
```

The contents of this file is not propagated by the cluster database; therefore, it should be installed on each node in the cluster along with the resource type's scripts. (If the

properties file is not installed on a given node and that node is used as the GUI server, the help text will not be displayed.)

Property Formats

In each resource type's properties file, you can have the following property:

```
resource_type.properFormat = introductory text
```

For each type-specific attribute, you can have the following properties:

- Label that will be displayed in the GUI:

```
resource_type.Attribute.label = GUI_label
```

- Help (glossary) text that will be linked to from each attribute's label:

```
resource_type.Attribute.glossary = glossary_key
```

- Information describing what the resource type is for and how it should be configured:

```
glossary_key = help text
```

Example Properties File

Following is an example properties file for the IP_Address resource type.

```
# IP_address

IP_address.properFormat = \
An IP address resource that belongs to a resource group can be \
used by clients to access the highly available resource group. \
As with any other type of resource, an IP address resource will \
be moved from one node to another when \
FailSafe detects a failure. \
The resource name for an IP address must follow standard dot \
notation. It should not be configured on any network interface. \
IP addresses that require name resolution are not valid IP_address \
resource names. \
For example, "192.0.2.22" could be the name of an IP_address \
resource.
```

```
IP_address.NetworkMask.label = \  
    Network Mask  
IP_address.NetworkMask.glossary = \  
    glossary.IP_address.NetworkMask  
glossary.IP_address.NetworkMask: \  
    <B>IP address network mask</B><P>\  
    The network mask of the IP address \  
    (for example, "0xffffffff00"). \  
    See the <B>ifconfig(1M)</B> reference page for more \  
    information.  
  
IP_address.interfaces.label = \  
    Interfaces  
IP_address.interfaces.glossary = \  
    glossary.IP_address.interfaces  
glossary.IP_address.interfaces: \  
    <B>IP address interfaces</B><P>\  
    A comma-separated list of interfaces on which the \  
    IP address can be configured \  
    (for example, "ec0,et0,ef0" or "hip0" or "lb0"). \  
    See the <B>ifconfig(1M)</B> reference page for more \  
    information.  
  
IP_address.BroadcastAddress.label = \  
    Broadcast Address  
IP_address.BroadcastAddress.glossary = \  
    glossary.IP_address.BroadcastAddress  
glossary.IP_address.BroadcastAddress: \  
    <B>IP address broadcast address</B><P>\  
    The broadcast address for the IP address \  
    (for example, "192.0.2.255"). \  
    See the <B>ifconfig(1M)</B> reference page for more \  
    information.
```

Testing a New Resource Type

After adding a new resource type, you should test it as follows:

1. Define a resource group that contains resources of the new type. Ensure that the group contains all of the resources on which the new resource type depends.
2. Bring the resource group online in the cluster using `cmgr` or the GUI.

For example, using `cmgr`:

```
cmgr> admin online resource_group new_rg in cluster TEST
```

3. Check the status of the resource group using `cmgr` or GUI after a few minutes.

For example:

```
cmgr> show status of resource_group new_rg in cluster TEST
```

4. If the resource group has been made online successfully, you will see output similar to the following:

```
State: Online
Error: No error
Owner: node1
```

5. If there are resource group errors, do the following:
 - Check the `srmd` logs (`/var/cluster/ha/log/srmd_nodename`) on the node on which the resource group is online.
 - Search for the string `ERROR` in the log file. There should be an error message about a resource in the resource group. The message also provides information about the action script that failed. For example:

```
Wed Nov 3 04:20:10.135 <E ha_srmd srm 12127:1 sa_process_tasks.c:627>
CI_FAILURE, ERROR: Action (exclusive) for resource (10.0.2.45) of type
(IP_address) failed with status (failed)
exclusive script failed for the resource 10.0.2.45 of resource type
IP_address. The status "failed"
indicates that the script returned an error.
```

- Check the script logs (`/var/cluster/ha/log/script_nodename` on the same node) for `IP_address` exclusive script errors.

- After the fixing the problems in the action script, perform an `offline_force` operation to clear the error. For example:

```
cmgr> admin offline_force resource_group new_rg in cluster TEST
```


Testing Scripts

This chapter describes how to test action scripts without running FailSafe. It also provides tips on how to debug problems that you may encounter. It covers the following:

- "General Testing and Debugging Techniques"
- "Debugging Notes" on page 80
- "Testing an Action Script" on page 81
- "Special Testing Considerations for the `monitor` Script" on page 83

Note: Parameters are passed to the action scripts as both input files and output files. Each line of the input file contains the resource name; the output file contains the resource name and the script exit status.

General Testing and Debugging Techniques

Some general testing and debugging techniques you can use during testing are as follows:

- To get debugging information, add the following line to each of your scripts in the main function of the script:

```
set -x
```

- To check that an application is running on a node:
 - Enter the following command on that node, where *application* is the name (or a portion of the name) of the executable for the application:

```
ps -ef | grep application
```

- Use appropriate commands provided by the application. For example, the FailSafe Informix option uses the Informix command `onstat`.
- To show the status of a resource, use the following `cmgr` command:

```
show status of resource resourcename of resource_type typename [in cluster clustername]
```

For example:

```
cmgr> show status of resource /hafs1/subdir of resource_type NFS in cluster nfs-cluster
```

```
State: Online  
Error: None  
Owner: hans2  
Flags: Resource is monitored locally
```

- To show the status of a node, use the following `cmgr` command:

```
show status of node nodename
```

For example:

```
cmgr> show status of node hans2
```

```
FailSafe status of node is UP.
```

```
Machine (hans2) is not configured for CXFS.
```

- To show the status of a resource group, use the following `cmgr` command:

```
show status of resource_group RG_name in cluster clustername
```

For example:

```
cmgr> show status of resource_group nfs-group1 in cluster nfs-cluster
```

```
State: Online  
Error: No error  
Owner: hans2
```

Debugging Notes

- The `exclusive` script returns an error when the resource is running in the local node. If the resource is actually running in the node, there is no `exclusive` action script bug.
- If the resource group does not become online on the primary node, it can be because of a `start` script error or a `monitor` script error on the primary node. The nature of the failure can be seen in the `srmd` logs of the primary node.

- If the action script failure status is `timeout`, resource type timeouts for the action should be increased. In the case of the `monitor` script, the check can be made more lightweight.
- The resource type action script timeouts are for a resource. So, if an action is performed on two resources, the script timeout is twice the configured resource type action timeout.
- If the resource group has a configuration error, check the `srmd` logs on the primary node for errors.
- The action scripts that use `${HA_LOG}` and `${HA_DBGLOG}` macros to log messages can find the messages in `/var/cluster/ha/log/script_nodename` file in each node in the cluster.

`HA_LOG` logs messages at log level 1 and `HA_DBGLOG` uses log level 11.

Testing an Action Script

To test an action script, do the following:

1. Create an input file, such as `/tmp/input`, that contains expected resource names. For example, to create a file that contains the resource named `disk1` do the following:

```
# echo "/disk1" > /tmp/input
```

2. Create an input parameter file, such as `/tmp/ipparamfile`, as follows:

```
# echo "ClusterName web-cluster" > /tmp/ipparamfile
```

3. Execute the action script as follows:

```
# ./start /tmp/input /tmp/output /tmp/ipparamfile
```

Note: The use of the input parameter file is optional.

4. Change the log level from `HA_NORMLVL` to `HA_DBGLVL` to allow messages written with `HA_DBGLOG` to be printed by adding the following line after the `set_global_variables` statement in your script:

```
HA_CURRENT_LOGLEVEL=$HA_DBGLVL
```

The output file will contain one of the following return values for the start, stop, monitor, and restart scripts:

```
HA_SUCCESS=0
HA_INVALID_ARGS=1
HA_CMD_FAILED=2
HA_NOTSUPPORTED=3
HA_NOCFGINFO=4
```

The output file will contain one of the following return values for the exclusive script:

```
HA_NOT_RUNNING=0
HA_RUNNING=2
```

Note: If you call the `exit_script` function prior to normal termination, it should be preceded by the `ha_write_status_for_resource` function and you should use the same return code that is logged to the output file.

Suppose you have a resource named `/disk1`. The syntax for the input and output files would be as follows:

- Input file: `<resourcename>`
- Output file: `<resourcename> <status>`

The following example shows:

- The exit status of the action script is 1
- The exit status of the resource is 2

Note: The use of `anonymous` indicates that the script was run manually. When the script is run by FailSafe, the full path to the script name is displayed.

```
# echo "/disk1" > /tmp/ipfile
# ./monitor /tmp/ipfile /tmp/opfile /tmp/ipparamfile
# echo $?
2
# cat /tmp/opfile
/disk1 2
# tail /var/cluster/ha/log/script_heb1
```

```
Tue Aug 25 11:32:57.437 <anonymous script 23787:0 Unknown:0> ./monitor:
./monitor called with /tmp/ipfile and /tmp/opfile
Tue Aug 25 11:32:58.118 <anonymous script 24556:0 Unknown:0> ./monitor:
check to see if /disk1 is mounted on /disk1
Tue Aug 25 11:32:58.433 <anonymous script 23811:0 Unknown:0> ./monitor:
/sbin/mount | grep /disk1 | grep /disk1 >> /dev/null 2>&l exited with
status 0
Tue Aug 25 11:32:58.665 <anonymous script 24124:0 Unknown:0> ./monitor:
stat mount point /disk1
Tue Aug 25 11:32:58.969 <anonymous script 23525:0 Unknown:0> ./monitor:
/sbin/stat /disk1 exited with status 0
Tue Aug 25 11:32:59.258 <anonymous script 24431:0 Unknown:0> ./monitor:
check the filesystem /disk1 is exported
Tue Aug 25 11:32:59.610 <anonymous script 6982:0 Unknown:0> ./monitor:
Tue Aug 25 11:32:59.917 <anonymous script 24040:0 Unknown:0> ./monitor:
awk '{print \$1}' /var/cluster/ha/tmp/exportfs.23762 | grep /disk1 exited
with status 1
Tue Aug 25 11:33:00.131 <anonymous script 24418:0 Unknown:0> ./monitor:
echo failed to find /disk1 in exported filesystem list:-
Tue Aug 25 11:33:00.340 <anonymous script 24236:0 Unknown:0> ./monitor:
echo /disk2
```

For additional information about a script's processing, see the */var/cluster/ha/log/script_nodename*.

Special Testing Considerations for the `monitor` Script

The `monitor` script tests the liveliness of applications and resources. The best way to test it is to induce a failure, run the script, and check if this failure is detected by the script; then repeat the process for another failure.

Use this checklist for testing a `monitor` script:

- Verify that the script detects failure of the application successfully
- Verify that the script always exits with a return value
- Verify that the script does not contain commands that can hang, such as using DNS for name resolution, or those that continue forever, such as `ping`

- Verify that the script completes before the time-out value specified in the configuration file
- Verify that the script's return codes are correct

During testing, measure the time it takes for a script to complete and adjust the monitoring times in your script accordingly. To get a good estimate of the time required for the script to execute, run it under different system load conditions.

Example: Requiring Confirmation Before Failover

Suppose you wanted to require that the system operator approve a failover before it takes place for an application that may not always recover automatically after a failure.

You could do this by creating a dummy resource type with the lowest order number possible: 1. This resource would be written to prompt the operator for confirmation before continuing with the failover. The operator would be prompted on any node in the cluster, depending upon where the resource is running. (The operator would also be prompted when the resource group is started for the first time, or when HA services are started after a reboot.)

A `start` script can wait indefinitely if you set its timeout value to a large-enough value. (The timeout value is specified in milliseconds.) If the start timeout expires, the resource group might go into error state. Therefore, the `start` script should wait for confirmation from the user only for certain amount of time. If the response is not received from the user, the `start` script can continue the failover or fail by writing `HA_CMD_FAILED` in the output file. If the `start` script fails, the user must manually recover the resource group.

Starting the FailSafe Manager GUI

There are several methods to start the GUI and connect to a node. For more information, see *FailSafe Administrator's Guide for SGI InfiniteStorage*.

Launch Methods

To start the GUI, use one of the following methods:

- On an IRIX system where the FailSafe GUI-client software (`sysadm_failsafe2.sw.client`) and desktop support software (`sysadm_failsafe.sw.desktop`) are installed, do one of the following:

Note: Do not use this method across a wide-area network (WAN) or virtual private network (VPN), or if the IRIX system has an R5000 or earlier CPU and less than 128-MB memory.

- Enter the following command line:

```
# /usr/sbin/fsmgr
```

(The `fsdetail` and `fstask` commands perform the identical function as `fsmgr`; these command names are kept for historical purposes.)

- Choose the following from the Toolchest:

```
System
  > FailSafe Manager
```

You must restart the Toolchest after installing FailSafe in order to see the **FailSafe** entry on the Toolchest display. Enter the following commands to restart the Toolchest:

```
# killall toolchest
# /usr/bin/X11/toolchest &
```

- On a PC or if you want to perform administration from a remote location via VPN or WAN, do the following:
 - Install a web server (such as Apache) and the `sysadm_failsafe2.sw.web` package on one of the administration nodes with an R5000 or later and at least 128-MB memory.
 - Install the Java2 v1.4.1 plug-in on your PC.
 - Close any existing Java windows and restart the Web browser on the PC.
 - Enter the following URL, where *server* is the name of a administration node in the pool:

`http://server/FailSafeManager/`
 - At the resulting webpage, click the FailSafe Manager icon.

Note: This method can be used on IRIX systems, but it is the preferred method only if you are using WAN or VPN. If you load the GUI using Netscape on IRIX and then switch to another page in Netscape, the FailSafe Manager GUI will not operate correctly. To avoid this problem, leave the FailSafe Manager GUI web page up and open a new Netscape window if you want to view another web page.

The following table describes the platforms where the GUI may be started, connected to, and displayed.

Table A-1 GUI Platforms

GUI Mode	Where You Start the GUI	Where You Connect the GUI	Where the GUI Displays
fsmgr(1) or Toolchest	An IRIX system (such as an SGI 2000 series, SGI O2 workstation, or Silicon Graphics Fuel visual workstation) with <code>sysadm_failsafe2.sw.client</code> and <code>sysadm_failsafe.sw.desktop</code> software installed	The FailSafe administration node in the pool that you want to use for cluster administration	The system where the GUI was invoked
Web	Any system with a web browser and Java 1.1 plug-in installed and enabled	The FailSafe administration node in the pool that you want to use for cluster administration	The same system with the web browser

Logging In

To ensure that the required GUI privileges are available for performing all of the tasks, you should log in to the GUI as `root`. However, some or all privileges can be granted to any other user using the GUI privilege tasks. (This functionality is also available with the Privilege Manager, part of the IRIX Interactive Desktop System Administration `sysadmdesktop` product. For more information, see the *Personal System Administration Guide*.)

A dialog box will appear, prompting you to log in to a FailSafe host. You can choose one of the following connection types:

- **Local** runs the server-side process on the local host instead of going over the network.
- **Direct** creates a direct socket connection using the `tcpmux` TCP protocol.
- **Remote Shell** connects to the server via a user-specified command shell, such as `rsh` or `ssh`.

Note: For a secure connection, choose **Remote Shell** and type a secure connection command using a utility such as `ssh`. Otherwise, the FailSafe Manager GUI will not encrypt communication and transferred passwords will be visible to users of the network.

- **Proxy** connects to the server through a firewall via a proxy server.

Making Changes from One Node

You should only make changes from one GUI process running at any given time; changes made by a second GUI process (a second invocation of `fsmgr`) may overwrite changes made by the first instance. However, multiple FailSafe Manager windows accessed via the **File** menu are all part of the same application process; you can make changes from any of these windows.

The FailSafe administration node to which you connect the GUI affects your view of the cluster. You should wait for a change to appear in the *details area* before making another change; the change is not guaranteed to be propagated across the cluster until it appears in the view area. The entire cluster status information is sent to every FailSafe administration node each time a change is made to the cluster database.

Using the Script Library

The purpose of the script library (`scriptlib`) is to simplify the FailSafe application interface so that users can use scripts and need not be aware of input and output file format.

The `/var/cluster/ha/common_scripts/scriptlib` file contains the library of environment variables (beginning with uppercase `HA_`) and functions (beginning with lowercase `ha_`) available for use in your action scripts.

Note: Do not change the contents of the `scriptlib` file.

This chapter describes functions that perform the following tasks, using samples from the `scriptlib` file:

- Set global definitions
- Check arguments
- Read an input file
- Execute a command
- Write status for a resource
- Get the value for a field
- Get resource information
- Print exclusivity check messages

File Formats

There are three file formats:

- *Input file*, which contains the list of resources that must be acted on by the executable; each resource must be specified on a separate line in the file.

- (Optional) *Output file*, in which the executable writes the return the status of each resource on a separate line, using the following format:

resource_name resource_status

There are corresponding lines for each line in the input file. The *resource_name* and *resource_status* fields are separated by whitespace. The resource status may be one of the following:

- HA_SUCCESS
- HA_RUNNING
- HA_NOT_RUNNING
- HA_INVALID_ARGS
- HA_CMD_FAILED
- HA_NOTSUPPORTED
- HA_NOCFGINFO (no configuration information)

If information about a resource is not present in the output file, SRMD assumes that the action on the resource has timed out. A nonzero value for the *resource_status* field is considered an error.

If the executable requires more information to perform the action on the resource, the information must be stored in the cluster database (CDB) in the local machine. The executables can use cluster database commands to extract information about the resource.

- *Input parameter file*, which contains the cluster name in the following format:

ClusterName *clustername*

Set Global Definitions

The `ha_set_global_defs()` function sets the global definitions for the environment variables shown in the following subsections.

The `HA_INFILE` and `HA_OUTFILE` variables set the input and output files for a script. These variables do not have global definitions, and are not set by the `ha_set_global_defs()` function.

Global Variable

`HA_HOSTNAME`

The output of the `uname` command with the `-n` option, which is the host name or node name. The node name is the name by which the system is known to communications networks.

Default: ``uname -n``

Command Location Variables

`HA_CMDSPATH`

Path to user commands.

Default: `/usr/cluster/bin`

`HA_PRIVCMDSPATH`

Path to privileged commands (those that can only be run by `root`).

Default: `/usr/sysadm/privbin`

`HA_LOGCMD`

Command used to log information.

Default: `ha_cilog`

HA_RESOURCEQUERYCMD

Resource query command. This is an internal command that is not meant for direct use in scripts; use the `ha_get_info()` function of `scriptlib` instead.

Default: `resourceQuery`

HA_SCRIPTTMPDIR

Location of the script temporary directory.

Default: `/tmp`

Database Location Variables

HA_CDB

Location of the cluster database.

Default: `/var/cluster/cdb/cdb.db`

Script Log Level Variables

HA_NORMLVL

Normal level of script logs.

Default: `0`

HA_DBGLVL

Debug level of script logs.

Default: `10`

Script Log Variables

HA_SCRIPTGROUP

Log for the script group.

Default: script

HA_SCRIPTSUBSYS

Log for the script subsystem.

Default: script

Script Logging Command Variables

HA_LOGQUERY_OUTPUT

Determine the current logging level for scripts.

Default:

```
`${HA_PRIVCMDSPATH}/loggroupQuery _NUM_LOG_GROUPS=1 \  
_LOG_GROUP_0=ha_script`
```

HA_DBGLOG

Command used to log debug messages from the scripts.

Default: ha_dbglog

HA_CURRENT_LOGLEVEL

Display the current log level. The default will be 0 (no script logging) if the loggroupQuery command fails or does not find configuration information.

Default: `echo \${HA_LOGQUERY_OUTPUT} | /usr/bin/awk '{print \$2}'`

HA_LOG

Command used to log the scripts.

Default: `ha_log`

Script Error Value Variables

HA_SUCCESS

Successful execution of the script. This variable is used by the `start`, `stop`, `restart`, and `monitor` scripts.

Default: 0

HA_NOT_RUNNING

The script is not running. This variable is used by `exclusive` scripts.

Default: 0

HA_INVALID_ARGS

An invalid argument was entered. This is used by all scripts.

Default: 1

HA_CMD_FAILED

A command called by the script has failed. This variable is used by the `start`, `stop`, `restart`, and `monitor` scripts.

Default: 2

HA_RUNNING

The script is running. This variable is used by `exclusive` scripts.

Default: 2

HA_NOTSUPPORTED

The specific action is not supported for this resource type. This is used by all scripts.

Default: 3

HA_NOCFGINFO

No configuration information was found. This is used by all scripts.

Default: 4

Check Arguments

An action script can have an input file (\$1 HA_INFILE), an output file (\$2 HA_OUTFILE), and a parameter file (\$3 HA_PARAMFILE). The parameter file is optional.

The `ha_check_args()` function checks the arguments specified for a script and sets the `$HA_INFILE` and `$HA_OUTFILE` variables accordingly.

If a parameter file exists, the `ha_check_args()` function reads the list of parameters from the file and sets the `$HA_CLUSTERNAME` variable.

In the following, long lines use the continuation character (`\`) for readability.

```
ha_check_args()
{
    ${HA_DBGLOG} "$HA_SCRIPTNAME called with $1, $2 and $3"

    if ! [ $# -eq 2 -o $# -eq 3 ]; then
        ${HA_LOG} "Incorrect number of arguments"
        return 1;
    fi

    if [ ! -r $1 ]; then
        ${HA_LOG} "file $1 is not readable or does not exist"
        return 1;
    fi
}
```

```
if [ ! -s $1 ]; then
    ${HA_LOG} "file $1 is empty"
    return 1;
fi
if [ $# -eq 3 ]; then
    HA_PARAMFILE=$3

    if [ ! -r $3 ]; then
        ${HA_LOG} "file $3 is not readable or does not exist"
        return 1;
    fi

    HA_CLUSTERNAME=`/usr/bin/awk '{ if ( $1 == "ClusterName" ) \
    print $2 }' ${HA_PARAMFILE}`
fi

HA_INFILE=$1
HA_OUTFILE=$2

return 0;
}
```

Read an Input File

The `ha_read_infile()` function reads the `$HA_INFILE` input file into the `$HA_RES_NAMES` variable, which specifies the list of resource names.

```
ha_read_infile()
{
    HA_RES_NAMES=" ";

    for HA_RESOURCE in `cat ${HA_INFILE}`
    do
        HA_TMP="${HA_RES_NAMES} ${HA_RESOURCE} ";
        HA_RES_NAMES=${HA_TMP};
    done
}
```

Execute a Command

The `ha_execute_cmd()` function executes the command specified by `$HA_CMD`, which is set in the action script. `$1` is the string to be logged. The function returns 1 on error and 0 on success. On errors, the standard output and standard error of the command is redirected to the log file.

```
ha_execute_cmd()
{
    OUTFILE=${HA_SCRIPTTMPDIR}/script.$$

    ${HA_DBGLOG} $1

    eval ${HA_CMD} > ${OUTFILE} 2>&1;

    ha_exit_code=$?;

    if [ $ha_exit_code -ne 0 ]; then
        ${HA_DBGLOG} `cat ${HA_SCRIPTTMPDIR}/script.$$`
    fi

    ${HA_DBGLOG} "${HA_CMD} exited with status $ha_exit_code";

    /sbin/rm ${OUTFILE}

    return $ha_exit_code;
}
```

The `ha_execute_cmd_ret()` function is similar to `ha_execute_cmd`, except that it places the command output in the location specified by `$HA_CMD_OUTPUT`.

```
ha_execute_cmd_ret()
{
    ${HA_DBGLOG} $1

    # REVISIT: Is it possible to redirect the output to a log
    HA_CMD_OUTPUT='${HA_CMD}';

    ha_exit_code=$?;

    ${HA_DBGLOG} "${HA_CMD} exited with status $ha_exit_code";

    return $ha_exit_code;
}
```

Write Status for a Resource

The `ha_write_status_for_resource()` function writes the status for a resource to the `$HA_OUTFILE` output file. `$1` is the resource name, and `$2` is the resource status.

```
ha_write_status_for_resource()
{
    echo $1 $2 >> $HA_OUTFILE;
}
```

Similarly, the `ha_write_status_for_all_resources()` function writes the status for all resources. `$HA_RES_NAMES` is the list of resource names.

```
ha_write_status_for_all_resources()
{
    for HA_RES in $HA_RES_NAMES
    do
        echo $HA_RES $1 >> $HA_OUTFILE;
    done
}
```

Get the Value for a Field

The `ha_get_field()` function obtains the field value from a string, where `$1` is the string and `$2` is the field name. The string format is as follows:

```
ha_get_field()
{
    HA_STR=$1
    HA_FIELD_NAME=$2
    ha_found=0;
    ha_field=1;

    for ha_i in $HA_STR
    do
        if [ $ha_field -eq 1 ]; then
            ha_field=0;
            if [ $ha_i = $HA_FIELD_NAME ]; then
                ha_found=1;
            fi
        else
            ha_field=1;
            if [ $ha_found -eq 1 ]; then
                HA_FIELD_VALUE=$ha_i
                return 0;
            fi
        fi
    done

    return 1;
}
```

Get the Value for Multiple Fields

The `ha_get_multi_fields()` function obtains the field values from a string, where `$1` is the string and `$2` is the field name. The string format is a series of name-value field pairs, where a name field is followed by the value of the name, separated by whitespace.

This function is typically used to extract dependency information. There may be multiple fields with the same name, so the string returned in `HA_FIELD_VALUE` may contain multiple values separated by white space. This appears as follows:

```
ha_get_multi_fields()
{
    HA_STR=$1
    HA_FIELD_NAME=$2
    ha_found=0;
    ha_field=1;

    for ha_i in $HA_STR
    do
        if [ $ha_field -eq 1 ]; then
            ha_field=0;
            if [ $ha_i = $HA_FIELD_NAME ]; then
                ha_found=1;
            fi
        else
            ha_field=1;
            if [ $ha_found -eq 1 ]; then
                if [ -z "$HA_FIELD_VALUE" ]; then
                    HA_FIELD_VALUE=$ha_i;
                else
                    HA_FIELD_VALUE="$HA_FIELD_VALUE $ha_i";
                fi;
            fi
            ha_found=0;
        fi
    done

    if [ -z "$HA_FIELD_VALUE" ]; then
        return 1;
    else
        return 0;
    fi
}
```

Get Resource Information

The `ha_get_info()` and `ha_get_info_debug()` functions read resource information. `$1` is the resource type, `$2` is the resource name, and `$3` is an optional parameter of any value that specifies a request for resource dependency information. Resource information is stored in the `HA_STRING` variable. If the `resourceQuery`

command fails, the HA_STRING is set to an invalid string, and future calls to `ha_get_info()` or `ha_get_info_debug()` return errors.

You can use `ha_get_info_debug()` while developing scripts.

```

ha_get_info()
{
    if [ -f /var/cluster/ha/resourceQuery.debug ]; then
        ha_get_info_debug $1 $2 $3
        return;
    fi

    if [ -n "$3" ]; then
        ha_doall="_ALL=true"
    else
        ha_doall=""
    fi

    # Retry resourceQuery command $SHA_RETRY_CMD_MAX times if $SHA_RETRY_CMD_ERR
    # is returned.
    ha_retry_count=1

    while [ $ha_retry_count -le $SHA_RETRY_CMD_MAX ];
    do
        if [ -n "${HA_CLUSTERNAME}" ]; then
            HA_STRING=`${HA_PRIVCMDSPATH}/${HA_RESOURCEQUERYCMD} \
                _CDB_DB=$HA_CDB _RESOURCE=$2 _RESOURCE_TYPE=$1 \
                $ha_doall _NO_LOGGING=true _CLUSTER=${HA_CLUSTERNAME}`
        else
            HA_STRING=`${HA_PRIVCMDSPATH}/${HA_RESOURCEQUERYCMD} \
                _CDB_DB=$HA_CDB _RESOURCE=$2 _RESOURCE_TYPE=$1 \
                $ha_doall _NO_LOGGING=true`
        fi

        ha_exit_code=$?

        if [ $ha_exit_code -ne 0 ]; then
            ${HA_LOG} "${HA_RESOURCEQUERYCMD}: resource name $2 resource type $1"
            ${HA_LOG} "Failed with error: ${HA_STRING}";
        fi

        if [ $ha_exit_code -ne $SHA_RETRY_CMD_ERR ]; then

```

```

        break;
    fi

    ha_retry_count=`expr $ha_retry_count + 1`

done

if [ -n "$ha_doall" ]; then
    echo $HA_STRING \
        | grep "No resource dependencies" > /dev/null 2>&1
    if [ $? -eq 0 ]; then
        HA_STRING=
    else
        HA_STRING=`echo $HA_STRING | /bin/sed -e "s/^. *Resource dependencies //"`
    fi
fi

return ${ha_exit_code};
}

```

The `ha_get_info` is a faster version of `ha_get_info_debug()`.

```

ha_get_info_debug()
{
    if [ -n "$3" ]; then
        ha_doall="_ALL=true"
    else
        ha_doall=""
    fi

    if [ -n "${HA_CLUSTERNAME}" ]; then
        HA_STRING=`${HA_PRIVCMDSPATH}/${HA_RESOURCEQUERYCMD} \
            _CDB_DB=$HA_CDB _RESOURCE=$2 _RESOURCE_TYPE=$1 \
            $ha_doall _CLUSTER=${HA_CLUSTERNAME}`
    else
        HA_STRING=`${HA_PRIVCMDSPATH}/${HA_RESOURCEQUERYCMD} \
            _CDB_DB=$HA_CDB _RESOURCE=$2 _RESOURCE_TYPE=$1 $ha_doall`
    fi
    ha_exit_code=$?

    if [ $? -ne 0 ]; then
        ${HA_LOG} "${HA_RESOURCEQUERYCMD}: resource name $2 resource type $1"
    fi
}

```

```

    ${HA_LOG} "Failed with error: ${HA_STRING}";
fi

if [ -n "$ha_doall" ]; then
    echo $HA_STRING \
        | grep "No resource dependencies" > /dev/null 2>&1
    if [ $? -eq 0 ]; then
        HA_STRING=
    else
        HA_STRING='echo $HA_STRING | /bin/sed -e "s/^.*Resource dependencies //"`
    fi
fi

return ${ha_exit_code};
}

```

Print Exclusivity Check Messages

The `ha_print_exclusive_status()` function prints exclusivity check messages to the log file. `$1` is the resource name and `$2` is the exit status.

```

ha_print_exclusive_status()
{
    if [ $? -eq $HA_NOT_RUNNING ]; then
        ${HA_LOG} "resource $1 exclusive status: NOT RUNNING"
    else
        ${HA_LOG} "resource $1 exclusive status: RUNNING"
    fi
}

```

The `ha_print_exclusive_status_all_resources()` function is similar, but it prints exclusivity check messages for all resources. `$HA_RES_NAMES` is the list of resource names.

```

ha_print_exclusive_status_all_resources()
{
    for HA_RES in $HA_RES_NAMES
    do
        ha_print_exclusive_status ${HA_RES} $1
    done
}

```

Glossary

action scripts

The set of scripts that determine how a resource is started, monitored, and stopped. There must be a set of action scripts specified for each resource type. The possible set of action scripts is: *exclusive*, *start*, *stop*, *monitor*, and *restart*.

active/backup configuration

A configuration in which all resource groups have the same primary node. The backup node does not run any highly available resource groups until a failover occurs.

cluster

A *cluster* is the set of systems (nodes) configured to work together as a single computing resource. A cluster is identified by a simple name and a cluster ID.

There is only one cluster that may be formed from a given pool of nodes.

Disks or logical units (LUNs) are assigned to clusters by recording the name of the cluster on the disk (or LUN). Thus, if any disk is accessible (via a Fibre Channel connection) from machines in multiple clusters, then those clusters must have unique names. When members of a cluster send messages to each other, they identify their cluster via the cluster ID. Thus, if two clusters will be sharing the same network for communications, then they must have unique cluster IDs.

Because of the above restrictions on cluster names and cluster IDs, and because cluster names and cluster IDs cannot be changed once the cluster is created (without deleting the cluster and recreating it), SGI advises that you choose unique names and cluster IDs for each of the clusters within your organization. Clusters that share a network and use XVM must have unique names.

cluster administration node

A node in a coexecution cluster that is installed with the `cluster_admin` software product, allowing the node to perform cluster administration tasks and contain a copy of the cluster database. Also known as a *CXFS administration node*.

cluster administrator

The person responsible for managing and maintaining a cluster.

cluster database

Contains configuration information about all resources, resource types, resource groups, failover policies, nodes, and the cluster.

cluster database membership

The group of nodes in the pool that are accessible to `fs2d` and therefore can receive cluster database updates; this may be a subset of the nodes defined in the pool. Also known as *user-space membership* and *fs2d membership*.

cluster ID

A unique number within your network in the range 1 through 128. The cluster ID is used by the IRIX kernel to make sure that it does not accept cluster information from any other cluster that may be on the network. The kernel does not use the database for communication, so it requires the cluster ID in order to verify cluster communications. This information in the kernel cannot be changed after it has been initialized; therefore, you must not change a cluster ID after the cluster has been defined.

cluster node

A node that is defined as part of the cluster. See also *node*.

cluster process group

A group of application instances in a distributed application that cooperate to provide a service.

For example, distributed lock manager instances in each node would form a process group. By forming a process group, they can obtain membership and reliable, ordered, atomic communication services. There is no relationship between a UNIX process group and a cluster process group.

collector host

The nodes in the FailSafe cluster itself from which you want to gather statistics, on which PCP for FailSafe has installed the collector agents.

control messages

Messages that cluster software sends between the nodes to request operations on or distribute information about nodes and resource groups. FailSafe sends control

messages for the purpose of ensuring that nodes and groups remain highly available. Control messages and heartbeat messages are sent through a node's network interfaces that have been attached to a control network. A node can be attached to multiple control networks.

control network

The network that connects nodes through their network interfaces (typically Ethernet) such that FailSafe can maintain a cluster's high availability by sending heartbeat messages and control messages through the network to the attached nodes. FailSafe uses the highest priority network interface on the control network; it uses a network interface with lower priority when all higher-priority network interfaces on the control network fail.

A node must have at least one control network interface for heartbeat messages and one for control messages (both heartbeat and control messages can be configured to use the same interface). A node can have no more than eight control network interfaces.

CXFS client administration node

A node that is installed with the `cluster_admin` software product, allowing the node to perform cluster administration tasks and contain a copy of the cluster database, but is not capable of coordinating CXFS metadata. FailSafe can run on a CXFS client-administration node.

CXFS client-only node

A node that is installed with the `cxfs_client.sw.base` software product; it does not run cluster administration daemons and is not capable of coordinating cluster activity and metadata. FailSafe cannot run on a client-only node.

CXFS server-capable administration node

A node in a coexecution cluster that is installed with the `cluster_admin` product and is also capable of coordinating CXFS metadata. FailSafe can run on a CXFS server-capable administration node.

database

See *cluster database*.

dependency list

See *resource dependency* or *resource type dependency*.

failover

The process of allocating a *resource group* to another *node* according to a *failover policy*. A failover may be triggered by the failure of a resource, a change in the FailSafe membership (such as when a node fails or starts), or a manual request by the administrator.

failover attribute

A string that affects the allocation of a resource group in a cluster. The administrator must specify system-defined attributes (such as `Auto_Failback` or `Controlled_Failback`), and can optionally supply site-specific attributes.

failover domain

The ordered list of nodes on which a particular *resource group* can be allocated. The nodes listed in the failover domain must be within the same cluster; however, the failover domain does not have to include every node in the cluster. The administrator defines the *initial failover domain* when creating a failover policy. This list is transformed into the *run-time failover domain* by the *failover script* the run-time failover domain is what is actually used to select the failover node. FailSafe stores the run-time failover domain and uses it as input to the next failover script invocation. The initial and run-time failover domains may be identical, depending upon the contents of the failover script. In general, FailSafe allocates a given resource group to the first node listed in the run-time failover domain that is also in the FailSafe membership; the point at which this allocation takes place is affected by the *failover attributes*.

failover policy

The method used by FailSafe to determine the destination node of a failover. A failover policy consists of a *failover domain*, *failover attributes*, and a *failover script*. A failover policy name must be unique within the *pool*.

failover script

A failover policy component that generates a *run-time failover domain* and returns it to the FailSafe process. The process applies the failover attributes and then selects the first node in the returned failover domain that is also in the current FailSafe membership.

FailSafe membership

The list of FailSafe nodes in a cluster on which FailSafe can make resource groups online. It differs from the CXFS membership. For more information about CXFS, see the *CXFS Administration Guide for SGI Infinite Storage*.

FailSafe database

See *cluster database*.

heartbeat messages

Messages that cluster software sends between the nodes that indicate a node is up and running. Heartbeat messages and *control messages* are sent through a node's network interfaces that have been attached to a control network. A node can be attached to multiple control networks.

heartbeat interval

Interval between heartbeat messages. The node timeout value must be at least 10 times the heartbeat interval for proper FailSafe operation (otherwise false failovers may be triggered). The higher the number of heartbeats (smaller heartbeat interval), the greater the potential for slowing down the network. Conversely, the fewer the number of heartbeats (larger heartbeat interval), the greater the potential for reducing availability of resources.

initial failover domain

The ordered list of nodes, defined by the administrator when a failover policy is first created, that is used the first time a cluster is booted. The ordered list specified by the initial failover domain is transformed into a *run-time failover domain* by the *failover script*; the run-time failover domain is used along with failover attributes to determine the node on which a resource group should reside. With each failure, the failover script takes the current run-time failover domain and potentially modifies it; the initial failover domain is never used again. Depending on the run-time conditions and contents of the failover script, the initial and run-time failover domains may be identical. See also *run-time failover domain*.

key/value attribute

A set of information that must be defined for a particular resource type. For example, for the resource type `filesystem` one key/value pair might be `mount_point=/fs1` where `mount_point` is the key and `fs1` is the value specific to the particular resource

being defined. Depending on the value, you specify either a `string` or `integer` data type. In the previous example, you would specify `string` as the data type for the value `fs1`.

log configuration

A log configuration has two parts: a *log level* and a *log file*, both associated with a *log group*. The cluster administrator can customize the location and amount of log output, and can specify a log configuration for all nodes or for only one node. For example, the `crsd` log group can be configured to log detailed level-10 messages to the `/var/cluster/ha/log/crsd-foo` log only on the node `foo` and to write only minimal level-1 messages to the `crsd` log on all other nodes.

log file

A file containing notifications for a particular *log group*. A log file is part of the *log configuration* for a log group. By default, log files reside in the `/var/cluster/ha/log` directory, but the cluster administrator can customize this. Note: FailSafe logs both normal operations and critical errors to `/var/adm/SYSLOG`, as well as to individual logs for specific log groups.

log group

A set of one or more FailSafe processes that use the same log configuration. A log group usually corresponds to one daemon, such as `gcd`.

log level

A number controlling the number of log messages that FailSafe will write into an associated log group's log file. A log level is part of the log configuration for a log group.

LUN

Logical unit number

monitor host

A workstation that has a display and is running the IRIS Desktop, on which PCP for FailSafe has installed the monitor client.

node

A *node* is an operating system (OS) image, usually an individual computer. (This use of the term *node* does not have the same meaning as a node in an SGI Origin 3000 or SGI 2000 system.)

A given node can be a member of only one pool (and therefore) only one cluster.

node ID

A 16-bit positive integer that uniquely defines a node. During node definition, FailSafe will assign a node ID if one has not been assigned by the cluster administrator. Once assigned, the node ID cannot be modified.

node timeout

If no heartbeat is received from a node in this period of time, the node is considered to be dead. The node timeout value must be at least 10 times the heartbeat interval for proper FailSafe operation (otherwise false failovers may be triggered).

notification command

The command used to notify the cluster administrator of changes or failures in the cluster, nodes, and resource groups. The command must exist on every node in the cluster.

offline resource group

A resource group that is not highly available in the cluster. To put a resource group in offline state, FailSafe stops the group (if needed) and stops monitoring the group. An offline resource group can be running on a node, yet not under FailSafe control. If the cluster administrator specifies the *detach only* option while taking the group offline, then FailSafe will not stop the group but will stop monitoring the group.

online resource group

A resource group that is highly available in the cluster. When FailSafe detects a failure that degrades the resource group availability, it moves the resource group to another node in the cluster. To put a resource group in online state, FailSafe starts the group (if needed) and begins monitoring the group. If the cluster administrator specifies the *attach only* option while bringing the group online, then FailSafe will not start the group but will begin monitoring the group.

owner host

A system that can control a node remotely, such as power-cycling the node. At run time, the owner host must be defined as a node in the pool.

owner TTY name

The device file name of the terminal port (TTY) on the *owner host* to which the system controller serial cable is connected. The other end of the cable connects to the node with the system controller port, so the node can be controlled remotely by the owner host.

plug-in

The set of software required to make an application highly available, including a resource type and action scripts. There are plug-ins provided with the base FailSafe release, optional plug-ins available for purchase from SGI, and customized plug-ins you can write using the instructions in this guide.

pool

The *pool* is the set of nodes from which a particular cluster may be formed. Only one cluster may be configured from a given pool, and it need not contain all of the available nodes. (Other pools may exist, but each is disjoint from the other. They share no node or cluster definitions.)

A pool is formed when you connect to a given node and define that node in the cluster database using the CXFS GUI or `cmgr` command. You can then add other nodes to the pool by defining them while still connected to the first node, or to any other node that is already in the pool. (If you were to connect to another node and then define it, you would be creating a second pool).

port password

The password for the system controller port, usually set once in firmware or by setting jumper wires. (This is not the same as the node's `root` password.)

powerfail mode

When powerfail mode is turned on, FailSafe tracks the response from a node's system controller as it makes reset requests to a node. When these requests fail to reset the node successfully, FailSafe uses heuristics to try to estimate whether the machine has been powered down. If the heuristic algorithm returns with success, FailSafe assumes

the remote machine has been reset successfully. When powerfail mode is turned off, the heuristics are not used and FailSafe may not be able to detect node power failures.

process membership

A list of process instances in a cluster that form a process group. There can multiple process groups per node.

re-MACing

The process of moving the physical medium access control (MAC) address of a network interface to another interface. It is done by using the `macconfig` command.

resource

A single physical or logical entity that provides a service to clients or other resources. For example, a resource can be a single disk volume, a particular network address, or an application such as a web server. A resource is generally available for use over time on two or more nodes in a cluster, although it can be allocated to only one node at any given time. Resources are identified by a resource name and a resource type. Dependent resources must be part of the same resource group and are identified in a resource dependency list.

resource dependency

The condition in which a resource requires the existence of other resources.

resource dependency list

A list of resources upon which a resource depends. Each resource instance must have resource dependencies that satisfy its resource type dependencies before it can be added to a resource group.

resource group

A collection of resources. A resource group is identified by a simple name; this name must be unique within a cluster. Resource groups cannot overlap; that is, two resource groups cannot contain the same resource. All interdependent resources must be part of the same resource group. If any individual resource in a resource group becomes unavailable for its intended use, then the entire resource group is considered unavailable. Therefore, a resource group is the unit of failover.

resource keys

Variables that define a resource of a given resource type. The action scripts use this information to start, stop, and monitor a resource of this resource type.

resource name

The simple name that identifies a specific instance of a resource type. A resource name must be unique within a given resource type.

resource type

A particular class of resource. All of the resources in a particular resource type can be handled in the same way for the purposes of failover. Every resource is an instance of exactly one resource type. A resource type is identified by a simple name; this name must be unique within a cluster. A resource type can be defined for a specific node or for an entire cluster. A resource type that is defined for a node overrides a cluster-wide resource type definition with the same name; this allows an individual node to override global settings from a cluster-wide resource type definition.

resource type dependency

A set of resource types upon which a resource type depends. For example, the `filesystem` resource type depends upon the `volume` resource type, and the `Netscape_web` resource type depends upon the `filesystem` and `IP_address` resource types.

resource type dependency list

A list of resource types upon which a resource type depends.

run-time failover domain

The ordered set of nodes on which the resource group can execute upon failures, as modified by the failover script. The run-time failover domain is used along with failover attributes to determine the node on which a resource group should reside. See also *initial failover domain*.

server-capable administration node

See *CXFS server-capable administration node*

start/stop order

Each resource type has a start/stop order, which is a nonnegative integer. In a resource group, the start/stop orders of the resource types determine the order in which the resources will be started when FailSafe brings the group online and will be stopped when FailSafe takes the group offline. The group's resources are started in increasing order, and stopped in decreasing order; resources of the same type are started and stopped in indeterminate order. For example, if resource type `volume` has order 10 and resource type `filesystem` has order 20, then when FailSafe brings a resource group online, all volume resources in the group will be started before all file system resources in the group.

system controller port

A port located on a node that provides a way to power-cycle the node remotely. Enabling or disabling a system controller port in the cluster database (CDB) tells FailSafe whether it can perform operations on the system controller port. (When the port is enabled, serial cables must attach the port to another node, the owner host.) System controller port information is optional for a node in the pool, but is required if the node will be added to a cluster; otherwise resources running on that node never will be highly available.

tie-breaker node

A node identified as a tie-breaker for FailSafe to use in the process of computing the FailSafe membership for the cluster, when exactly half the nodes in the cluster are up and can communicate with each other. If a tie-breaker node is not specified, FailSafe will use the node with the lowest node ID in the cluster as the tie-breaker node.

type-specific attribute

Required information used to define a resource of a particular resource type. For example, for a resource of type `filesystem` you must enter attributes for the resource's volume name (where the file system is located) and specify options for how to mount the file system (for example, as readable and writable).

Index

A

- action scripts
 - definition of the term, 7
 - examples, 24
 - failure of, 12
 - format
 - basic action, 21
 - completion, 23
 - exit status, 20
 - header, 18
 - overview, 18
 - read input file, 22
 - read resource information, 19, 20
 - set global variables, 22
 - set local variables, 19
 - verify arguments, 22
 - monitoring
 - frequency, 16
 - necessity of, 15
 - testing examples, 17
 - types, 16
 - preparation for writing scripts, 14
 - resource types provided, 14
 - set of scripts, 7
 - successful execution results, 11
 - templates, 14
 - testing, 81
 - writing steps, 23
- administrative commands, 5
- agents, 34
- Auto_Failback failover attribute, 40
- Auto_Recovery failover attribute, 40

C

- check arguments, 97
- cluster database security, 59
- cluster_mgr command (see cmgr), 66
- cmgr command, 66
- cmond process configuration, 34
- command execution, 99
- command path, 93
- commands, 5
- communicate with the network interface agent
 - daemon, 6
- configurations
 - N+1, 50
 - N+2, 52
 - N+M, 53
- Controlled_Failback failover attribute, 40
- Critical_RG failover attribute, 41
- CXFS resource type, 2

D

- database location, 94
- database security, 59
- debug script messages, 95
- debugging information in action scripts, 79
- DMF resource type, 2
- domain, 38

E

- environment variables, 93
- exclusive script
 - definition, 7
 - example, 31

execute a command, 99
exit status in action scripts, 20
exit_script(), 20, 82
exit_status value, 20

F

failover attributes, 39
failover domain, 38
failover policy
 contents, 37
 examples
 N+1, 50
 N+2, 51
 N+M, 53
 failover attributes, 39
 failover domain, 38
failover script, 41
 description, 41
 interface, 48
field value, 101
file locking and unlocking, 6
File menu, 90
filesystem resource type, 2
fsdetail, 87
fsdetail (fsmgr), 90
fsmgr, 87, 90
fstask, 87
fstask (fsmgr), 90

G

get_xxx_info(), 20
global definition setting, 93
global variables, 22
GUI
 multiple instances, 90
 starting, 89

H

HA_CDB, 94
ha_check_args(), 22, 97
ha_cilog command, 6
HA_CMD_FAILED, 96
HA_CMDSPATH, 93
HA_CURRENT_LOGLEVEL, 95
HA_DBGLOG, 95
HA_DBGLVL, 94
ha_exec2, 12
ha_exec2 command, 6
ha_execute_cmd(), 99
ha_execute_cmd_ret(), 100
ha_execute_lock command, 6
ha_filelock command, 6
ha_fileunlock command, 6
ha_get_field(), 101
ha_get_info(), 20, 102
ha_get_multi_fields(), 20
HA_HOSTNAME, 93
ha_http_ping2 command, 6
ha_ifdadmin command, 6
HA_INVALID_ARGS, 96
HA_LOG, 96
HA_LOGCMD, 93
HA_LOGQUERY_OUTPUT, 95
ha_macconfig2 command, 6
HA_NOCFGINFO, 97
HA_NORMLVL, 94
HA_NOT_RUNNING, 96
HA_NOTSUPPORTED, 97
ha_print_exclusive_status(), 105
ha_print_exclusive_status_all_resources(), 105
HA_PRIVCMDSPATH, 93
ha_read_infile(), 22, 98
HA_RESOURCEQUERYCMD, 94
HA_RUNNING, 96
HA_SCRIPTGROUP, 95
HA_SCRIPTSUBSYS, 95
HA_SCRIPTTMPDIR, 94

HA_SUCCESS, 96
ha_write_status_for_all_resources(), 100
ha_write_status_for_resource, 21
ha_write_status_for_resource(), 100
high availability characteristics, 3
hostname, 93

I

Informix resource type, 2
initial failover domain, 38
InPlace_Recovery failover attribute, 40
input file, 98
IP address resource type, 2

L

lock a file, 6
log messages, 6
logs, 95

M

MAC_address resource type, 2
message logging, 6
monitor script
 definition, 7
 example, 28
monitoring
 agents, 34
 failure, 16
 frequency, 17
 necessity of, 15
 processes, 6
 script testing, 83
 testing examples, 17
 types, 15

N

Netscape node check, 6
Netscape resource type, 2
NFS resource type, 2
node status, 80
Node_Failures_Only failover attribute, 41
nodename output, 93

O

Oracle resource type, 2
order ranges for resource types, 57
ordered failover script, 41
overview of the programming steps, 2

P

path to user commands, 93
plug-ins, 1
print exclusivity check messages, 105
Privilege Manager, 89
privileged command path, 93
process monitoring, 6
programming steps overview, 2
properties file, 73

R

read an input file, 98
resource group states, 11
resource information
 obtaining, 102
 read into an action script, 20
resource query command, 94
resource type
 cmgr use, 66
 GUI use, 61

- information for a new resource type, 56
- order ranges, 57
- restart mode, 57
- script templates, 72
- script use, 71
- restart mode, 57
- restart script
 - definition, 8
 - example, 33
- root command path, 93
- run-time failover domain, 38

S

- Samba resource type, 2
- script group log, 95
- script library, 91
- script testing
 - action scripts, 81
 - monitoring script considerations, 83
 - techniques, 79
- script.\$\$ suffix, 23
- scriptlib file, 91
- scripts.
 - See "action scripts or failover script", 18
- security of the cluster database, 59
- set_global_variables(), 22
- set_local_variables() section of an action script, 19
- start script
 - definition, 7
 - example, 24
- std_unlimited resource type, 2
- status of a node, 80
- stop script
 - definition, 7
 - example, 26
- sysadm_failsafe2.sw.desktop, 87
- sysadmdesktop, 89

T

- templates
 - action scripts, 14
 - resource type script definition, 72
- testing scripts
 - See "script testing", 79
- TMF resource type, 2
- Toolchest, 87

U

- uname, 93
- unlock a file, 6
- user command path, 93
- user privileges, 59

V

- value for a field, 101
 - /var/cluster/cmgr-templates/
 - cmgr-create-resource_type directory, 72
 - /var/cluster/cmon/process_groups directory, 34
 - /var/cluster/ha/
 - resource_types directory, 58
 - resource_types/<resource_type>/<resource_type>, 73
 - /var/cluster/ha/policies directory, 41
- volume resource type, 2

W

- write status for a resource, 100

X

- XFS resource type, 2
- XLV resource type, 2

XVM resource type, 2