# **Red Hat Enterprise Linux 5**

## **Online Storage Reconfiguration Guide**

For Red Hat Enterprise Linux 5



Red Hat Enterprise Linux Engineering Mike Christie Tom Coughlan Don Domingo Jacquelynn East Rob Evers Development Community Pasi Kärkkäinen

#### Abstract

This document outlines the different procedures involved in reconfiguring storage devices while the system is running. The document currently addresses iSCSI and Fibre Channel interconnects. Other interconnects may be added in future versions of this document.

1. Introduction	2
1.1. Document Conventions	
1.2. Getting Help and Giving Feedback	6
2. Fibre Channel	
2.1. Fibre Channel API	6
2.2. Native Fibre Channel Drivers and Capabilities	7
3. iSCSI	8
3.1. iSCSI API	8
4. Persistent Naming	9
4.1. WWID	9
4.2. UUID and Other Persistent Identifiers	11
5. Removing a Storage Device	11
6. Removing a Path to a Storage Device	
7. Adding a Storage Device or Path	13
8. Configuring a Fibre-Channel Over Ethernet Interface	14
9. Scanning Storage Interconnects	15
10. iSCSI Discovery Configuration	16

11.	Configuring iSCSI Offload and Interface Binding	17
	11.1. Viewing Available iface Configurations	17
	11.2. Configuring an iface for Software iSCSI	19
	11.3. Configuring an iface for iSCSI Offload	19
	11.4. Binding/Unbinding an iface to a Portal	20
12.	Scanning iSCSI Interconnects	
13.	Logging In to an iSCSI Target	23
14.	Resizing an Online Logical Unit	23
	14.1. Resizing Fibre Channel Logical Units	
	14.2. Resizing an iSCSI Logical Unit	24
	14.3. Updating the Size of Your Multipath Device	25
15.	Adding/Removing a Logical Unit Through rescan-scsi-bus.sh	26
16.	Modifying Link Loss Behavior	26
	16.1. Fibre Channel	26
	16.2. iSCSI Settings With dm-multipath	27
	16.3. iSCSI Root	30
	Controlling the SCSI Command Timer and Device Status	
18.	Troubleshooting	32
A. I	Revision History	33

#### Index

## **1. Introduction**

It is often desirable to add, remove or re-size storage devices while the operating system is running, and without rebooting. This manual outlines the procedures that may be used to reconfigure storage devices on Red Hat Enterprise Linux 5 host systems while the system is running. It covers iSCSI and Fibre Channel storage interconnects; other interconnect types may be added it the future.

33

The *Online Storage Reconfiguration guide* focuses on adding, removing, modifying, and monitoring storage devices. It does not discuss the Fibre Channel or iSCSI protocols in detail. For more information about these protocols, refer to other documentation.

This manual assumes that you have advanced working knowledge of Red Hat Enterprise Linux 5, along with first-hand experience in managing storage devices in Linux. Before consulting this book, verify if your host bus adapter vendor or hardware vendor have their own documentation. It is recommended that you consult such documents in conjunction with this manual.

#### Note

This manual makes reference to various **sysfs** objects. Red Hat advises that the **sysfs** object names and directory structure are subject to change in major Red Hat Enterprise Linux 5 releases. This is because the upstream Linux kernel does not provide a stable internal API. For guidelines on how to reference **sysfs** objects in a transportable way, refer to the document **Documentation/sysfs-rules.txt** in the kernel source tree for guidelines.

## Warning

Online storage reconfiguration must be done carefully. System failures or interruptions during the process can lead to unexpected results. Red Hat advises that you reduce system load to the maximum extent possible during the change operations. This will reduce the chance of I/O errors, out-of-memory errors, or similar errors occurring in the midst of a configuration change. The following sections provide more specific guidelines regarding this.

In addition, Red Hat recommends that you back up all data before performing online storage reconfiguration.

## **1.1. Document Conventions**

This manual uses several conventions to highlight certain words and phrases and draw attention to specific pieces of information.

In PDF and paper editions, this manual uses typefaces drawn from the *Liberation Fonts*<sup>1</sup> set. The Liberation Fonts set is also used in HTML editions if the set is installed on your system. If not, alternative but equivalent typefaces are displayed. Note: Red Hat Enterprise Linux 5 and later includes the Liberation Fonts set by default.

#### 1.1.1. Typographic Conventions

Four typographic conventions are used to call attention to specific words and phrases. These conventions, and the circumstances they apply to, are as follows.

#### Mono-spaced Bold

Used to highlight system input, including shell commands, file names and paths. Also used to highlight keycaps and key combinations. For example:

To see the contents of the file **my\_next\_bestselling\_novel** in your current working directory, enter the **cat my\_next\_bestselling\_novel** command at the shell prompt and press **Enter** to execute the command.

The above includes a file name, a shell command and a keycap, all presented in mono-spaced bold and all distinguishable thanks to context.

Key combinations can be distinguished from keycaps by the hyphen connecting each part of a key combination. For example:

Press Enter to execute the command.

Press **Ctrl+Alt+F2** to switch to the first virtual terminal. Press **Ctrl+Alt+F1** to return to your X-Windows session.

The first paragraph highlights the particular keycap to press. The second highlights two key combinations (each a set of three keycaps with each set pressed simultaneously).

<sup>&</sup>lt;sup>1</sup> https://fedorahosted.org/liberation-fonts/

If source code is discussed, class names, methods, functions, variable names and returned values mentioned within a paragraph will be presented as above, in **mono-spaced bold**. For example:

File-related classes include **filesystem** for file systems, **file** for files, and **dir** for directories. Each class has its own associated set of permissions.

#### **Proportional Bold**

This denotes words or phrases encountered on a system, including application names; dialog box text; labeled buttons; check-box and radio button labels; menu titles and sub-menu titles. For example:

Choose System  $\rightarrow$  Preferences  $\rightarrow$  Mouse from the main menu bar to launch Mouse Preferences. In the Buttons tab, click the Left-handed mouse check box and click Close to switch the primary mouse button from the left to the right (making the mouse suitable for use in the left hand).

To insert a special character into a **gedit** file, choose **Applications**  $\rightarrow$  **Accessories** 

 $\rightarrow$  Character Map from the main menu bar. Next, choose Search  $\rightarrow$  Find... from the Character Map menu bar, type the name of the character in the Search field and click Next. The character you sought will be highlighted in the Character Table. Double-click this highlighted character to place it in the Text to copy field and then click the

Copy button. Now switch back to your document and choose  $\textbf{Edit} \rightarrow \textbf{Paste}$  from the **gedit** menu bar.

The above text includes application names; system-wide menu names and items; application-specific menu names; and buttons and text found within a GUI interface, all presented in proportional bold and all distinguishable by context.

#### Mono-spaced Bold Italic or Proportional Bold Italic

Whether mono-spaced bold or proportional bold, the addition of italics indicates replaceable or variable text. Italics denotes text you do not input literally or displayed text that changes depending on circumstance. For example:

To connect to a remote machine using ssh, type **ssh** *username@domain.name* at a shell prompt. If the remote machine is **example.com** and your username on that machine is john, type **ssh** john@example.com.

The **mount** -o **remount** *file-system* command remounts the named file system. For example, to remount the **/home** file system, the command is **mount** -o **remount /home**.

To see the version of a currently installed package, use the **rpm** -**q** *package* command. It will return a result as follows: *package-version-release*.

Note the words in bold italics above — username, domain.name, file-system, package, version and release. Each word is a placeholder, either for text you enter when issuing a command or for text displayed by the system.

Aside from standard usage for presenting the title of a work, italics denotes the first use of a new and important term. For example:

Publican is a DocBook publishing system.

#### 1.1.2. Pull-quote Conventions

Terminal output and source code listings are set off visually from the surrounding text.

Output sent to a terminal is set in mono-spaced roman and presented thus:

booksDesktopdocumentationdraftsmssphotosstuffsvnbooks\_testsDesktop1downloadsimagesnotesscriptssvgs

Source-code listings are also set in **mono-spaced roman** but add syntax highlighting as follows:

```
package org.jboss.book.jca.ex1;
import javax.naming.InitialContext;
public class ExClient
{
   public static void main(String args[])
      throws Exception
   {
     InitialContext iniCtx = new InitialContext();
     Object
                  ref = iniCtx.lookup("EchoBean");
     EchoHome
                    home = (EchoHome) ref;
     Echo
                    echo
                          = home.create();
     System.out.println("Created Echo");
     System.out.println("Echo.echo('Hello') = " + echo.echo("Hello"));
  }
}
```

#### 1.1.3. Notes and Warnings

Finally, we use three visual styles to draw attention to information that might otherwise be overlooked.

## Note

Notes are tips, shortcuts or alternative approaches to the task at hand. Ignoring a note should have no negative consequences, but you might miss out on a trick that makes your life easier.

## Important

Important boxes detail things that are easily missed: configuration changes that only apply to the current session, or services that need restarting before an update will apply. Ignoring a box labeled 'Important' will not cause data loss but may cause irritation and frustration.



Warnings should not be ignored. Ignoring warnings will most likely cause data loss.

## **1.2. Getting Help and Giving Feedback**

#### 1.2.1. Do You Need Help?

If you experience difficulty with a procedure described in this documentation, visit the Red Hat Customer Portal at *http://access.redhat.com*. Through the customer portal, you can:

- search or browse through a knowledgebase of technical support articles about Red Hat products.
- submit a support case to Red Hat Global Support Services (GSS).
- access other product documentation.

Red Hat also hosts a large number of electronic mailing lists for discussion of Red Hat software and technology. You can find a list of publicly available mailing lists at *https://www.redhat.com/mailman/listinfo*. Click on the name of any mailing list to subscribe to that list or to access the list archives.

#### 1.2.2. We Need Feedback!

If you find a typographical error in this manual, or if you have thought of a way to make this manual better, we would love to hear from you! Please submit a report in Bugzilla: *http://bugzilla.redhat.com/* against the product **Red\_Hat\_Enterprise\_Linux.** 

When submitting a bug report, be sure to mention the manual's identifier: *Online\_Storage\_Reconfiguration\_Guide* 

If you have a suggestion for improving the documentation, try to be as specific as possible when describing it. If you have found an error, please include the section number and some of the surrounding text so we can find it easily.

## 2. Fibre Channel

This section discusses the Fibre Channel API, native Red Hat Enterprise Linux 5 Fibre Channel drivers, and the Fibre Channel capabilities of these drivers.

#### 2.1. Fibre Channel API

Below is a list of /sys/class/ directories that contain files used to provide the userspace API. In each item, host numbers are designated by *H*, bus numbers are *B*, targets are *T*, logical unit numbers (LUNs) are *L*, and remote port numbers are *R*.



If your system is using multipath software, Red Hat recommends that you consult your hardware vendor before changing any of the values described in this section.

Transport: /sys/class/fc\_transport/targetH:B:T/

- **port\_id** 24-bit port ID/address
- node\_name 64-bit node name
- **port\_name** 64-bit port name

Remote Port: /sys/class/fc\_remote\_ports/rport-H:B-R/

- port\_id
- node\_name
- port\_name
- •

**dev\_loss\_tmo** — number of seconds to wait before marking a link as "bad". Once a link is marked bad, IO running on its corresponding path (along with any new IO on that path) will be failed.

The default dev\_loss\_tmo value varies, depending on which driver/device is used. If a Qlogic adapter is used, the default is 35 seconds, while if an Emulex adapter is used, it is 30 seconds. The dev\_loss\_tmo value can be changed via the scsi\_transport\_fc module parameter dev\_loss\_tmo, although the driver can override this timeout value.

The maximum **dev\_loss\_tmo** value is 600 seconds. If **dev\_loss\_tmo** is set to zero or any value greater than 600, the driver's internal timeouts will be used instead.

• fast io fai

**fast\_io\_fail\_tmo** — length of time to wait before failing IO executed when a link problem is detected. IO that reaches the driver will fail. If IO is in a blocked queue, it will not be failed until **dev\_loss\_tmo** expires and the queue is unblocked.

#### Host: /sys/class/fc\_host/hostH/

- port\_id
- •

**issue\_lip** — instructs the driver to rediscover remote ports.

#### 2.2. Native Fibre Channel Drivers and Capabilities

Red Hat Enterprise Linux 5 ships with the following native fibre channel drivers:

- lpfc
- qla2xxx
- zfcp
- mptfc

*Table 1, "Fibre-Channel API Capabilities"* describes the different fibre-channel API capabilities of each native Red Hat Enterprise Linux 5 driver. X denotes support for the capability.

	lpfc	qla2xxx	zfcp	mptfc
Transport <b>port_id</b>	Х	X	X	X
Transport node_name	X	X	X	X
Transport <b>port_name</b>	X	X	X	Х
Remote Port dev_loss_tmo	X	X	X	Х
Remote Port fast_io_fail_t	X mo	X <sup>1</sup>		
Host port_id	Х	X	X	X
Host issue_lip	Х	X		

#### Table 1. Fibre-Channel API Capabilities

<sup>1</sup> Supported as of Red Hat Enterprise Linux 5.4

## 3. iSCSI

This section describes the iSCSI API and the **iscsiadm** utility. Before using the **iscsiadm** utility, install the **iscsi-initiator-utils** package first; to do so, run **yum install iscsi-initiator-utils**.

In addition, the iSCSI service must be running in order to discover or log in to targets. To start the iSCSI service, run **service iscsi start** 

#### 3.1. iSCSI API

To get information about running sessions, run:

```
iscsiadm -m session -P 3
```

This command displays the session/device state, session ID (sid), some negotiated parameters, and the SCSI devices accessible through the session.

For shorter output (for example, to display only the sid-to-node mapping), run:

```
iscsiadm -m session -P 0
```

or

```
iscsiadm -m session
```

These commands print the list of running sessions with the format:

driver [sid] target\_ip:port,target\_portal\_group\_tag proper\_target\_name

For example:

```
iscsiadm -m session
```

tcp [2] 10.15.84.19:3260,2 iqn.1992-08.com.netapp:sn.33615311 tcp [3] 10.15.85.19:3260,3 iqn.1992-08.com.netapp:sn.33615311

For more information about the iSCSI API, refer to /usr/share/doc/iscsi-initiatorutils-version/README.

## 4. Persistent Naming

The operating system issues I/O to a storage device by referencing the path that is used to reach it. For SCSI devices, the path consists of the following:

- · PCI identifier of the host bus adapter (HBA)
- channel number on that HBA
- the remote SCSI target address
- the Logical Unit Number (LUN)

This path-based address is not persistent. It may change any time the system is reconfigured (either by on-line reconfiguration, as described in this manual, or when the system is shutdown, reconfigured, and rebooted). It is even possible for the path identifiers to change when no physical reconfiguration has been done, as a result of timing variations during the discovery process when the system boots, or when a bus is re-scanned.

The operating system provides several non-persistent names to represent these access paths to storage devices. One is the **/dev/sd** name; another is the **major:minor** number. A third is a symlink maintained in the **/dev/disk/by-path/** directory. This symlink maps from the path identifier to the current **/dev/sd** name. For example, for a Fibre Channel device, the PCI info and **Host:BusTarget:LUN** info may appear as follows:

pci-0000:02:0e.0-scsi-0:0:0:0 -> ../../sda

For iSCSI devices, **by-path**/ names map from the target name and portal information to the **sd** name.

It is generally *not* appropriate for applications to use these path-based names. This is because the storage device these paths reference may change, potentially causing incorrect data to be written to the device. Path-based names are also not appropriate for multipath devices, because the path-based names may be mistaken for separate storage devices, leading to uncoordinated access and unintended modifications of the data.

In addition, path-based names are system-specific. This can cause unintended data changes when the device is accessed by multiple systems, such as in a cluster.

For these reasons, several persistent, system-independent, methods for identifying devices have been developed. The following sections discuss these in detail.

#### 4.1. WWID

The *World Wide Identifier* (WWID) can be used in reliably identifying devices. It is a persistent, system-independent ID that the SCSI Standard requires from all SCSI devices. The WWID identifier is

guaranteed to be unique for every storage device, and independent of the path that is used to access the device.

This identifier can be obtained by issuing a SCSI Inquiry to retrieve the *Device Identification Vital Product Data* (page **0x83**) or *Unit Serial Number* (page **0x80**). The mappings from these WWIDs to the current /dev/sd names can be seen in the symlinks maintained in the /dev/disk/by-id/ directory.

For example, a device with a page **0x83** identifier would have:

```
scsi-3600508b400105e210000900000490000 -> ../../sda
```

Or, a device with a page **0x80** identifier would have:

```
scsi-SSEAGATE_ST373453LW_3HW1RHM6 -> ../../sda
```

Red Hat Enterprise Linux 5 automatically maintains the proper mapping from the WWID-based device name to a current /dev/sd name on that system. Applications can use the /dev/disk/by-id/ name to reference the data on the disk, even if the path to the device changes, and even when accessing the device from different systems.

If there are multiple paths from a system to a device, **device-mapper-multipath** uses the WWID to detect this. **Device-mapper-multipath** then presents a single "pseudo-device" in /dev/mapper/ wwid, such as /dev/mapper/3600508b400105df70000e00000ac0000.

The command **multipath** -1 shows the mapping to the non-persistent identifiers: *Host:Channel:Target:LUN*, /dev/sd name, and the **major:minor** number.

```
3600508b400105df70000e0000ac0000 dm-2 vendor,product
[size=20G][features=1 queue_if_no_path][hwhandler=0][rw]
\_ round-robin 0 [prio=0][active]
\_ 5:0:1:1 sdc 8:32 [active][undef]
\_ 6:0:1:1 sdg 8:96 [active][undef]
\_ round-robin 0 [prio=0][enabled]
\_ 5:0:0:1 sdb 8:16 [active][undef]
\_ 6:0:0:1 sdf 8:80 [active][undef]
```

**Device-mapper-multipath** automatically maintains the proper mapping of each WWID-based device name to its corresponding **/dev/sd** name on the system. These names are persistent across path changes, and they are consistent when accessing the device from different systems.

When the **user\_friendly\_names** feature (of **device-mapper-multipath**) is used, the WWID is mapped to a name of the form **/dev/mapper/mpathn**. By default, this mapping is maintained in the file **/var/lib/multipath/bindings**. These **mpathn** names are persistent as long as that file is maintained.



The multipath bindings file (by default, **/var/lib/multipath/bindings**) must be available at boot time. If **/var** is a separate filesystem from **/**, then you must change the default location of the file. For more information, refer to *http://kbase.redhat.com/faq/docs/DOC-17650*.



If you use **user\_friendly\_names**, then additional steps are required to obtain consistent names in a cluster. Refer to the *Consistent Multipath Device Names*<sup>2</sup> section in the *Using Device-Mapper Multipath*<sup>3</sup> book.

In addition to these persistent names provided by the system, you can also use **udev** rules to implement persistent names of your own, mapped to the WWID of the storage. For more information about this, refer to *http://kbase.redhat.com/faq/docs/DOC-7319*.

## 4.2. UUID and Other Persistent Identifiers

If a storage device contains a filesystem, then that filesystem may provide one or both of the following:

- Universally Unique Identifier (UUID)
- · Filesystem label

These identifiers are persistent, and based on metadata written on the device by certain applications. They may also be used to access the device using the symlinks maintained by the operating system in the /dev/disk/by-label/ (e.g. boot -> ../../sda1 ) and /dev/disk/by-uuid/ (e.g. f8bf09e3-4c16-4d91-bd5e-6f62da165c08 -> ../../sda1) directories.

**md** and LVM write metadata on the storage device, and read that data when they scan devices. In each case, the metadata contains a UUID, so that the device can be identified regardless of the path (or system) used to access it. As a result, the device names presented by these facilities are persistent, as long as the metadata remains unchanged.

## 5. Removing a Storage Device

Before removing access to the storage device itself, it is advisable to back up data from the device first. Afterwards, flush I/O and remove all operating system references to the device (as described below). If the device uses multipathing, then do this for the multipath "pseudo device" (*Section 4.1, "WWID"*) and each of the identifiers that represent a path to the device. If you are only removing a path to a multipath device, and other paths will remain, then the procedure is simpler, as described in *Section 7, "Adding a Storage Device or Path"*.

Removal of a storage device is not recommended when the system is under memory pressure, since the I/O flush will add to the load. To determine the level of memory pressure, run the command **vmstat 1 100**; device removal is not recommended if:

• Free memory is less than 5% of the total memory in more than 10 samples per 100 (the command **free** can also be used to display the total memory).

 <sup>&</sup>lt;sup>2</sup> http://www.redhat.com/docs/en-US/Red\_Hat\_Enterprise\_Linux/5.4/html/DM\_Multipath/multipath\_consistent\_names.html
 <sup>3</sup> http://www.redhat.com/docs/en-US/Red\_Hat\_Enterprise\_Linux/5.4/html/DM\_Multipath/index.html

• Swapping is active (non-zero si and so columns in the vmstat output).

The general procedure for removing all access to a device is as follows:

#### Procedure 1. Ensuring a Clean Device Removal

- 1. Close all users of the device and backup device data as needed.
- 2. Use **umount** to unmount any file systems that mounted the device.
- 3. Remove the device from any **md** and LVM volume using it. If the device is a member of an LVM Volume group, then it may be necessary to move data off the device using the **pvmove** command, then use the **vgreduce** command to remove the physical volume, and (optionally) **pvremove** to remove the LVM metadata from the disk.
- 4. If the device uses multipathing, run **multipath** -1 and note all the paths to the device. Afterwards, remove the multipathed device using **multipath** -f device.
- 5. Run **blockdev** –**flushbufs** *device* to flush any outstanding I/O to all paths to the device. This is particularly important for raw devices, where there is no **umount** or **vgreduce** operation to cause an I/O flush.
- 6. Remove any reference to the device's path-based name, like /dev/sd, /dev/disk/by-path or the major:minor number, in applications, scripts, or utilities on the system. This is important in ensuring that different devices added in the future will not be mistaken for the current device.
- Finally, remove each path to the device from the SCSI subsystem. To do so, use the command echo 1 > /sys/block/device-name/device/delete where device-name may be sde, for example.

Another variation of this operation is **echo 1** > /sys/class/scsi\_device/h:c:t:1/ device/delete, where *h* is the HBA number, *c* is the channel on the HBA, *t* is the SCSI target ID, and *1* is the LUN.



You can determine the *device-name*, HBA number, HBA channel, SCSI target ID and LUN for a device from various commands, such as **lsscsi**, **scsi\_id**, **multipath** -1, and **ls** -1 /dev/ disk/by-\*.

After performing *Procedure 1, "Ensuring a Clean Device Removal"*, a device can be physically removed safely from a running system. It is not necessary to stop I/O to other devices while doing so.

Other procedures, such as the physical removal of the device, followed by a rescan of the SCSI bus (as described in *Section 9, "Scanning Storage Interconnects"*) to cause the operating system state to be updated to reflect the change, are not recommended. This will cause delays due to I/O timeouts, and devices may be removed unexpectedly. If it is necessary to perform a rescan of an interconnect, it must be done while I/O is paused, as described in *Section 9, "Scanning Storage Interconnects"*.

## 6. Removing a Path to a Storage Device

If you are removing a path to a device that uses multipathing (without affecting other paths to the device), then the general procedure is as follows:

#### Procedure 2. Removing a Path to a Storage Device

- 1. Remove any reference to the device's path-based name, like **/dev/sd** or **/dev/disk/by-path** or the **major:minor** number, in applications, scripts, or utilities on the system. This is important in ensuring that different devices added in the future will not be mistaken for the current device.
- 2. Take the path offline using echo offline > /sys/block/sda/device/state.

This will cause any subsequent IO sent to the device on this path to be failed immediately. **Device-mapper-multipath** will continue to use the remaining paths to the device.

 Remove the path from the SCSI subsystem. To do so, use the command echo 1 > /sys/ block/device-name/device/delete where device-name may be sde, for example (as described in Procedure 1, "Ensuring a Clean Device Removal").

After performing *Procedure 2, "Removing a Path to a Storage Device"*, the path can be safely removed from the running system. It is not necessary to stop I/O while this is done, as **device-mapper-multipath** will re-route I/O to remaining paths according to the configured path grouping and failover policies.

Other procedures, such as the physical removal of the cable, followed by a rescan of the SCSI bus to cause the operating system state to be updated to reflect the change, are not recommended. This will cause delays due to I/O timeouts, and devices may be removed unexpectedly. If it is necessary to perform a rescan of an interconnect, it must be done while I/O is paused, as described in *Section 9*, *"Scanning Storage Interconnects"*.

## 7. Adding a Storage Device or Path

When adding a device, be aware that the path-based device name (/dev/sd name, major:minor number, and /dev/disk/by-path name, for example) the system assigns to the new device may have been previously in use by a device that has since been removed. As such, ensure that all old references to the path-based device name have been removed. Otherwise, the new device may be mistaken for the old device.

The first step in adding a storage device or path is to physically enable access to the new storage device, or a new path to an existing device. This is done using vendor-specific commands at the Fibre Channel or iSCSI storage server. When doing so, note the LUN value for the new storage that will be presented to your host. If the storage server is Fibre Channel, also take note of the *World Wide Node Name* (WWNN) of the storage server, and determine whether there is a single WWNN for all ports on the storage server. If this is not the case, note the *World Wide Port Name* (WWPN) for each port that will be used to access the new LUN.

Next, make the operating system aware of the new storage device, or path to an existing device. The recommended command to use is:

#### echo "c t 1" > /sys/class/scsi\_host/hosth/scan

In the previous command, h is the HBA number, c is the channel on the HBA, t is the SCSI target ID, and 1 is the LUN.

The older form of this command, echo "scsi add-single-device 0 0 0 0" > /proc/ scsi/scsi, is deprecated.

For Fibre Channel storage servers that implement a single WWNN for all ports, you can determine the correct *h*,*c*,and *t* values (i.e. HBA number, HBA channel, and SCSI target ID) by searching for the WWNN in **sysfs**. For example, if the WWNN of the storage server is **0x5006016090203181**, use:

#### grep 5006016090203181 /sys/class/fc\_transport/\*/node\_name

This should display output similar to the following:

/sys/class/fc\_transport/target5:0:2/node\_name:0x5006016090203181
/sys/class/fc\_transport/target5:0:3/node\_name:0x5006016090203181
/sys/class/fc\_transport/target6:0:2/node\_name:0x5006016090203181
/sys/class/fc\_transport/target6:0:3/node\_name:0x5006016090203181

This indicates there are four Fibre Channel routes to this target (two single-channel HBAs, each leading to two storage ports). Assuming a LUN value is **56**, then the following command will configure the first path:

#### echo "0 2 56" > /sys/class/scsi\_host/host5/scan

This must be done for each path to the new device.

For Fibre Channel storage servers that do not implement a single WWNN for all ports, you can determine the correct HBA number, HBA channel, and SCSI target ID by searching for each of the WWPNs in **sysfs**.

Another way to determine the HBA number, HBA channel, and SCSI target ID is to refer to another device that is already configured on the same path as the new device. This can be done with various commands, such as **lsscsi**, **scsi\_id**, **multipath -1**, and **ls -1 /dev/disk/by-\***. This information, plus the LUN number of the new device, can be used as shown above to probe and configure that path to the new device.

After adding all the SCSI paths to the device, execute the **multipath** command, and check to see that the device has been properly configured. At this point, the device can be added to **md**, LVM, **mkfs**, or **mount**, for example.

If the steps above are followed, then a device can safely be added to a running system. It is not necessary to stop I/O to other devices while this is done. Other procedures involving a rescan (or a reset) of the SCSI bus, which cause the operating system to update its state to reflect the current device connectivity, are not recommended while storage I/O is in progress.

## 8. Configuring a Fibre-Channel Over Ethernet Interface

Setting up and deploying a Fibre-channel over ethernet (FCoE) interface requires two packages:

- · fcoe-utils
- dcbd

Once these packages are installed, perform the following procedure to enable FCoE over a virtual LAN (VLAN):

Procedure 3. Configuring an ethernet interface to use FCoE

1. Configure a new VLAN (101) by creating a new network script for it. The easiest way to do so is to copy the network script of an ethernet interface (eth3) to a new one with the **101** file name suffix, as in:

cp /etc/sysconfig/network-scripts/ifcfg-eth3 /etc/sysconfig/networkscripts/ifcfg-eth3.101

2. Open **/etc/sysconfig/network-scripts/ifcfg-eth3.101**. Edit it to ensure that the following settings are configured accordingly:

DEVICE=eth3.101 VLAN=yes ONBOOT=yes

3. Start the data center bridging daemon (dcbd) using the following command:

#### /etc/init.d/dcbd start

4. Use the **dcbtool** utility to enable data center bridging and FCoE on the ethernet interface using the following commands:

dcbtool sc eth3 dcb on

dcbtool sc eth3 app:fcoe e:1

These commands will only work if no other changes have been made to the **dcbd** settings for the ethernet interface.

5. Start FCoE using the command **/etc/init.d/fcoe start**. The fibre-channel device should appear shortly, assuming all other settings on the fabric are correct.

After correctly configuring the ethernet interface to use FCoE, Red Hat recommends that set FCoE and **dcbd** to run at startup. To do so, use **chkconfig**, as in:

chkconfig dcbd on

chkconfig fcoe on

## 9. Scanning Storage Interconnects

There are several commands available that allow you to reset and/or scan one or more interconnects, potentially adding and removing multiple devices in one operation. This type of scan can be disruptive, as it can cause delays while I/O operations timeout, and remove devices unexpectedly. As such, Red Hat recommends that this type of scan be used *only when necessary*. In addition, the following restrictions must be observed when scanning storage interconnects:

- 1. All I/O on the effected interconnects must be paused and flushed before executing the procedure, and the results of the scan checked before I/O is resumed.
- As with removing a device, interconnect scanning is not recommended when the system is under memory pressure. To determine the level of memory pressure, run the command vmstat 1 100; interconnect scanning is not recommended if free memory is less than 5% of the total memory in

more than 10 samples per 100. It is also not recommended if swapping is active (non-zero **si** and **so** columns in the **vmstat** output). The command **free** can also display the total memory.

The following commands can be used to scan storage interconnects.

#### echo "1" > /sys/class/fc\_host/host/issue\_lip

This operation performs a *Loop Initialization Protocol (LIP)* and then scans the interconnect and causes the SCSI layer to be updated to reflect the devices currently on the bus. A LIP is, essentially, a bus reset, and will cause device addition and removal. This procedure is necessary to configure a new SCSI target on a Fibre Channel interconnect.

Bear in mind that **issue\_lip** is an asynchronous operation. The command may complete before the entire scan has completed. You must monitor **/var/log/messages** to determine when it is done.

The **lpfc** and **qla2xxx** drivers support **issue\_lip**. For more information about the API capabilities supported by each driver in Red Hat Enterprise Linux, refer to *Table 1, "Fibre-Channel API Capabilities"*.

#### /usr/bin/rescan-scsi-bus.sh

This script is included in Red Hat Enterprise Linux 5.4 and all future updates. By default, this script scans all the SCSI buses on the system, updating the SCSI layer to reflect new devices on the bus. The script provides additional options to allow device removal and the issuing of LIPs. For more information about this script (including known issues), refer to *Section 15, "Adding/Removing a Logical Unit Through rescan-scsi-bus.sh"*.

#### echo "- - -" > /sys/class/scsi\_host/hosth/scan

This is the same command described in *Section 7*, *"Adding a Storage Device or Path"* to add a storage device or path. In this case, however, the channel number, SCSI target ID, and LUN values are replaced by wildcards. Any combination of identifiers and wildcards is allowed, allowing you to make the command as specific or broad as needed. This procedure will add LUNs, but not remove them.

#### rmmod driver-name or modprobe driver-name

These commands completely re-initialize the state of all interconnects controlled by the driver. Although this is extreme, it may be appropriate in some situations. This may be used, for example, to re-start the driver with a different module parameter value.

## **10. iSCSI Discovery Configuration**

The default iSCSI configuration file is **/etc/iscsi/iscsid.conf**. This file contains iSCSI settings used by **iscsid** and **iscsiadm**.

During target discovery, the **iscsiadm** tool uses the settings in **/etc/iscsi/iscsid.conf** to create two types of records:

#### Node records in /var/lib/iscsi/nodes

When logging into a target, **iscsiadm** uses the settings in this file.

#### Discovery records in /var/lib/iscsi/discovery\_type

When performing discovery to the same destination, iscsiadm uses the settings in this file.

Before using different settings for discovery, delete the current discovery records (i.e. /var/lib/ iscsi/discovery\_type) first. To do this, use the following command:

```
iscsiadm -m discovery -t discovery_type -p target_IP:port -o delete<sup>4</sup>
```

Here, *discovery\_type* can be either sendtargets, isns, or fw.<sup>5</sup>

There are two ways to reconfigure discovery record settings:

• Edit the /etc/iscsi/iscsid.conf file directly prior to performing a discovery. Discovery settings use the prefix discovery; to view them, run:

```
iscsiadm -m discovery -t discovery_type -p target_IP:port
```

• Alternatively, **iscsiadm** can also be used to directly change discovery record settings, as in:

```
iscsiadm -m discovery -t discovery_type -p target_IP:port -o update -n
setting -v %value
```

Refer to man iscsiadm for more information on available settings and valid values for each.

After configuring discovery settings, any subsequent attempts to discover new targets will use the new settings. Refer to *Section 12, " Scanning iSCSI Interconnects"* for details on how to scan for new iSCSI targets.

For more information on configuring iSCSI target discovery, refer to the **man** pages of **iscsiadm** and **iscsid**. The **/etc/iscsi/iscsid.conf** file also contains examples on proper configuration syntax.

## 11. Configuring iSCSI Offload and Interface Binding

This chapter describes how to set up iSCSI interfaces in order to bind a session to a NIC port when using software iSCSI. It also describes how to set up interfaces for use with network devices that support offloading; namely, devices from Chelsio, Broadcom and ServerEngines.

The network subsystem can be configured to determine the path/NIC that iSCSI interfaces should use for binding. For example, if portals and NICs are set up on different subnets, then it is not necessary to manually configure iSCSI interfaces for binding.

Before attempting to configure an iSCSI interface for binding, run the following command first:

#### ping -I ethX target\_IP<sup>4</sup>

If **ping** fails, then you will not be able to bind a session to a NIC. If this is the case, check the network settings first.

#### **11.1.** Viewing Available iface Configurations

Red Hat Enterprise Linux 5.5 supports iSCSI offload and interface binding for the following iSCSI initiator implementations:

Software iSCSI — like the scsi\_tcp and ib\_iser modules, this stack allocates an iSCSI host instance (i.e. scsi\_host) per session, with a single connection per session. As a result, /sys/ class\_scsi\_host and /proc/scsi will report a scsi\_host for each connection/session you are logged into.

<sup>&</sup>lt;sup>4</sup> The *target\_IP* and *port* variables refer to the IP address and port combination of a target/portal, respectively. For more information, refer to *Section 3.1, "iSCSI API"* and *Section 12, " Scanning iSCSI Interconnects*".

<sup>&</sup>lt;sup>5</sup> For details on different types of discovery, refer to the *DISCOVERY TYPES* section of **man iscsiadm**.

Offload iSCSI — like the Chelsio cxgb3i, Broadcom bnx2i and ServerEngines be2iscsi modules, this stack allocates a scsi\_host for each PCI device. As such, each port on a host bus adapter will show up as a different PCI device, with a different scsi\_host per HBA port.

To manage both types of initiator implementations, **iscsiadm** uses the **iface** structure. With this structure, an **iface** configuration must be entered in **/var/lib/iscsi/ifaces** for each HBA port, software iSCSI, or network device (**ethX**) used to bind sessions.

To view available **iface** configurations, run **iscsiadm** -**m iface**. This will display **iface** information in the following format:

iface\_name transport\_name, hardware\_address, ip\_address, net\_ifacename, initiator\_name

Refer to the following table for an explanation of each value/setting.

Table	2.	iface	Settings
Table	<u> </u>	naoc	Coungo

Setting	Description
iface_name	iface configuration name.
transport_name	Name of driver
hardware_address	MAC address
ip_address	IP address to use for this port
net_iface_name	Name used for the <b>vlan</b> or alias binding of a software iSCSI session. For iSCSI offloads, <b>net_iface_name</b> will be <b><empty></empty></b> because this value is not persistent across reboots.
initiator_name	This setting is used to override a default name for the initiator, which is defined in <b>/etc/</b> <b>iscsi/initiatorname.iscsi</b>

The following is a sample output of the iscsiadm -m iface command:

iface0 qla4xxx,00:c0:dd:08:63:e8,20.15.0.7,default,iqn.2005-06.com.redhat:madmax iface1 qla4xxx,00:c0:dd:08:63:ea,20.15.0.9,default,iqn.2005-06.com.redhat:madmax

For software iSCSI, each **iface** configuration must have a unique name (with less than 65 characters). The **iface\_name** for network devices that support offloading appears in the format *transport\_name.hardware\_name*.

For example, the sample output of **iscsiadm** -**m iface** on a system using a Chelsio network card might appear as:

```
default tcp,<empty>,<empty>,<empty>,<empty>,
iser iser,<empty>,<empty>,<empty>,<empty>,
cxgb3i.00:07:43:05:97:07 cxgb3i,00:07:43:05:97:07,<empty>,<empty>,<empty>,
```

It is also possible to display the settings of a specific **iface** configuration in a more friendly way. To do so, use the option **-I** *iface\_name*. This will display the settings in the following format:

```
iface.setting = value
```

Using the previous example, the **iface** settings of the same Chelsio video card (i.e. **iscsiadm** -m **iface** -**I** cxgb3i.00:07:43:05:97:07) would appear as:

```
# BEGIN RECORD 2.0-871
iface.iscsi_ifacename = cxgb3i.00:07:43:05:97:07
iface.net_ifacename = <empty>
iface.ipaddress = <empty>
iface.hwaddress = 00:07:43:05:97:07
iface.transport_name = cxgb3i
iface.initiatorname = <empty>
# END RECORD
```

#### 11.2. Configuring an iface for Software iSCSI

As mentioned earlier, an **iface** configuration is required for each network object that will be used to bind a session.

Before

To create an **iface** configuration for software iSCSI, run the following command:

```
iscsiadm -m iface -I iface_name --op=new
```

This will create a new *empty* **iface** configuration with a specified *iface\_name*. If an existing **iface** configuration already has the same *iface\_name*, then it will be overwritten with a new, empty one.

To configure a specific setting of an **iface** configuration, use the following command:

```
iscsiadm -m iface -I iface_name --op=update -n iface.setting -v hw_address
```

For example, to set the MAC address (hardware\_address) of iface0 to 00:0F:1F:92:6B:BF, run:

```
iscsiadm -m iface -I iface0 - -op=update -n iface.hwaddress -v
00:0F:1F:92:6B:BF
```

Warning

Do not use **default** or **iser** as **iface** names. Both strings are special values used by **iscsiadm** for backward compatibility. Any manually-created **iface** configurations named **default** or **iser** will disable backwards compatibility.

#### **11.3.** Configuring an iface for iSCSI Offload

By default, **iscsiadm** will create an **iface** configuration for each Chelsio, Broadcom, and ServerEngines port. To view available **iface** configurations, use the same command for doing so in software iSCSI, i.e. **iscsiadm** -m **iface**.

Before using the **iface** of a network card for iSCSI offload, first set the IP address (*target\_IP*<sup>4</sup>) that the device should use. For ServerEngines devices that use the **be2iscsi** driver (i.e. ServerEngines iSCSI HBAs), the IP address is configured in the ServerEngines BIOS setup screen.

For Chelsio and Broadcom devices, the procedure for configuring the IP address is the same as for any other **iface** setting. So to configure the IP address of the **iface**, use:

#### iscsiadm -m iface -I iface\_name -o update -n iface.ipaddress -v target\_IP

For example, to set the **iface** IP address of a Chelsio card (with **iface** name **cxgb3i.00:07:43:05:97:07**) to **20.15.0.66**, use:

iscsiadm -m iface -I cxgb3i.00:07:43:05:97:07 -o update -n iface.ipaddress -v 20.15.0.66

#### 11.4. Binding/Unbinding an iface to a Portal

Whenever **iscsiadm** is used to scan for interconnects, it will first check the **iface.transport** settings of each **iface** configuration in **/var/lib/iscsi/ifaces**. The **iscsiadm** utility will then bind discovered portals to any **iface** whose **iface.transport** is **tcp**.

This behavior was implemented for compatibility reasons. To override this, use the **-I** *iface\_name* to specify which portal to bind to an *iface*, as in:

iscsiadm -m discovery -t st -p target\_IP:port -I iface\_name -P 1<sup>4</sup>

By default, the **iscsiadm** utility will not automatically bind any portals to **iface** configurations that use offloading. This is because such **iface** configurations will not have **iface.transport** set to **tcp**. As such, the **iface** configurations of Chelsio, Broadcom, and ServerEngines ports need to be manually bound to discovered portals.

It is also possible to prevent a portal from binding to any existing **iface**. To do so, use **default** as the **iface\_name**, as in:

iscsiadm -m discovery -t st -p IP:port -I default -P 1

To remove the binding between a target and **iface**, use:

```
iscsiadm -m node -targetname proper_target_name -I iface0 --op=delete<sup>6</sup>
```

To delete all bindings for a specific **iface**, use:

```
iscsiadm -m node -I iface_name --op=delete
```

To delete bindings for a specific portal (e.g. for EqualLogic targets), use:

iscsiadm -m node -p IP:port -I iface\_name --op=delete

Note

If there are no **iface** configurations defined in **/var/lib/iscsi/iface** and the **-I** option is not used, **iscsiadm** will allow the network subsystem to decide which device a specific portal should use.

## 12. Scanning iSCSI Interconnects

For iSCSI, if the targets send an iSCSI async event indicating new storage is added, then the scan is done automatically. Cisco MDS<sup>TM</sup> and EMC Celerra<sup>TM</sup> support this feature.

However, if the targets do not send an iSCSI async event, you need to manually scan them using the **iscsiadm** utility. Before doing so, however, you need to first retrieve the proper **--targetname** and the **--portal** values. If your device model supports only a single logical unit and portal per target, use **iscsiadm** to issue a **sendtargets** command to the host, as in:

<sup>&</sup>lt;sup>6</sup> Refer to Section 12, "Scanning iSCSI Interconnects" for information on proper\_target\_name.

#### iscsiadm -m discovery -t sendtargets -p target\_IP:port<sup>4</sup>

The output will appear in the following format:

```
target_IP:port,target_portal_group_tag proper_target_name
```

For example, on a target with a *proper\_target\_name* of iqn.1992-08.com.netapp:sn.33615311 and a *target\_IP:port* of 10.15.85.19:3260, the output may appear as:

```
10.15.84.19:3260,2 iqn.1992-08.com.netapp:sn.33615311
10.15.85.19:3260,3 iqn.1992-08.com.netapp:sn.33615311
```

In this example, the target has two portals, each using *target\_ip:ports* of 10.15.84.19:3260 and 10.15.85.19:3260.

To see which **iface** configuration will be used for each session, add the **-P 1** option. This option will print also session information in tree format, as in:

```
Target: proper_target_name
    Portal: target_IP:port,target_portal_group_tag
    Iface Name: iface_name
```

For example, with iscsiadm -m discovery -t sendtargets -p 10.15.85.19:3260 -P 1, the output may appear as:

```
Target: iqn.1992-08.com.netapp:sn.33615311
Portal: 10.15.84.19:3260,2
Iface Name: iface2
Portal: 10.15.85.19:3260,3
Iface Name: iface2
```

This means that the target iqn.1992-08.com.netapp:sn.33615311 will use iface2 as its iface configuration.

With some device models (e.g. from EMC and Netapp), however, a single target may have multiple logical units and/or portals. In this case, issue a **sendtargets** command to the host first to find new portals on the target. Then, rescan the existing sessions using:

#### iscsiadm -m session --rescan

You can also rescan a specific session by specifying the session's SID value, as in:

#### iscsiadm -m session -r SID --rescan<sup>7</sup>

If your device supports multiple targets, you will need to issue a **sendtargets** command to the hosts to find new portals for *each* target. Then, rescan existing sessions to discover new logical units on existing sessions (i.e. using the **--rescan** option).

<sup>&</sup>lt;sup>7</sup> For information on how to retrieve a session's SID value, refer to Section 3.1, "iSCSI API".

Important

The **sendtargets** command used to retrieve **--targetname** and **--portal** values overwrites the contents of the **/var/lib/iscsi/nodes** database. This database will then be repopulated using the settings in **/etc/iscsi/iscsid.conf**. However, this will not occur if a session is currently logged in and in use.

To safely add new targets/portals or delete old ones, use the **-o new** or **-o delete** options, respectively. For example, to add new targets/portals without overwriting **/var/lib/iscsi/nodes**, use the following command:

```
iscsiadm -m discovery -t st -p target_IP -o new
```

To delete /var/lib/iscsi/nodes entries that the target did not display during discovery, use:

iscsiadm -m discovery -t st -p target\_IP -o delete

You can also perform both tasks simultaneously, as in:

```
iscsiadm -m discovery -t st -p target_IP -o delete -o new
```

The sendtargets command will yield the following output:

ip:port,target\_portal\_group\_tag proper\_target\_name

For example, given a device with a single target, logical unit, and portal, with **equallogic-iscsi1** as your **target\_name**, the output should appear similar to the following:

10.16.41.155:3260,0 iqn.2001-05.com.equallogic:6-8a0900-ac3fe0101-63aff113e344a4a2-dl585-03-1

Note that *proper\_target\_name* and *ip:port,target\_portal\_group\_tag* are identical to the values of the same name in *Section 3.1, "iSCSI API*".

At this point, you now have the proper **--targetname** and **--portal** values needed to manually scan for iSCSI devices. To do so, run the following command:

## iscsiadm --mode node --targetname proper\_target\_name --portal ip:port,target\_portal\_group\_tag\ --login <sup>8</sup>

Using our previous example (where *proper\_target\_name* is **equallogic-iscsi1**), the full command would be:

<sup>&</sup>lt;sup>8</sup> This is a single command split into multiple lines, to accommodate printed and PDF versions of this document. All concatenated lines — preceded by the backslash (\) — should be treated as one command, sans backslashes.

```
iscsiadm --mode node --targetname \ iqn.2001-05.com.equallogic:6-8a0900-
ac3fe0101-63aff113e344a4a2-dl585-03-1 \ --portal 10.16.41.155:3260,0 --
login<sup>8</sup>
```

## 13. Logging In to an iSCSI Target

As mentioned in Section 3, "iSCSI", the iSCSI service must be running in order to discover or log into targets. To start the iSCSI service, run:

```
service iscsi start
```

When this command is executed, the iSCSI **init** scripts will automatically log into targets where the **node.startup** setting is configured as **automatic**. This is the default value of **node.startup** for all targets.

To prevent automatic login to a target, set **node.startup** to **manual**. To do this, run the following command:

```
iscsiadm -m node --targetname proper_target_name -p target_IP:port -o
update -n node.startup -v manual
```

Deleting the entire record will also prevent automatic login. To do this, run:

```
iscsiadm -m node --targetname proper_target_name -p target_IP:port -o
delete
```

To automatically mount a file system from an iSCSI device on the network, add a partition entry for the mount in **/etc/fstab** with the **\_netdev** option. For example, to automatically mount the iSCSI device **sdb** to **/mount/iscsi** during startup, add the following line to **/etc/fstab**:

/dev/sdb /mnt/iscsi ext3 \_netdev 0 0

To manually log in to an iSCSI target, use the following command:

```
iscsiadm -m node --targetname proper_target_name -p target_IP:port -1
```

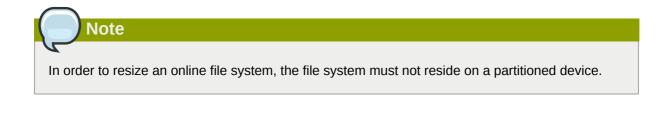
## Note

The *proper\_target\_name* and *target\_IP:port* refer to the full name and IP address/port combination of a target. For more information, refer to *Section 3.1, "iSCSI API"* and *Section 12, "Scanning iSCSI Interconnects"*.

## 14. Resizing an Online Logical Unit

In most cases, fully resizing an online *logical unit* involves two things: resizing the logical unit itself and reflecting the size change in the corresponding multipath device (if multipathing is enabled on the system).

To resize the online logical unit, start by modifying the logical unit size through the array management interface of your storage device. This procedure differs with each array; as such, consult your storage array vendor documentation for more information on this.



#### 14.1. Resizing Fibre Channel Logical Units

After modifying the online logical unit size, re-scan the logical unit to ensure that the system detects the updated size. To do this for Fibre Channel logical units, use the following command:

echo 1 > /sys/block/sdX/device/rescan

Important

To re-scan fibre channel logical units on a system that uses multipathing, execute the aforementioned command for each sd device (i.e. **sd1**, **sd2**, and so on) that represents a path for the multipathed logical unit. To determine which devices are paths for a multipath logical unit, use **multipath -11**; then, find the entry that matches the logical unit being resized. It is advisable that you refer to the WWID of each entry to make it easier to find which one matches the logical unit being resized.

## 14.2. Resizing an iSCSI Logical Unit

After modifying the online logical unit size, re-scan the logical unit to ensure that the system detects the updated size. To do this for iSCSI devices, use the following command:

```
iscsiadm -m node --targetname target_name -R<sup>4</sup>
```

Replace *target\_name* with the name of the target where the device is located.

Note

You can also re-scan iSCSI logical units using the following command:

iscsiadm -m node -R -I interface

Replace *interface* with the corresponding interface name of the resized logical unit (for example, **iface0**). This command performs two operations:

- It scans for new devices in the same way that the command echo "- -" > /sys/ class/scsi\_host/host/scan does (refer to Section 12, " Scanning iSCSI Interconnects").
- It re-scans for new/modified logical units the same way that the command echo 1 > /sys/ block/sdX/device/rescan does. Note that this command is the same one used for rescanning fibre-channel logical units.

## 14.3. Updating the Size of Your Multipath Device

If multipathing is enabled on your system, you will also need to reflect the change in logical unit size to the logical unit's corresponding multipath device (*after* resizing the logical unit). For Red Hat Enterprise Linux 5.3 (or later), you can do this through **multipathd**. To do so, first ensure that **multipathd** is running using **service multipathd status**. Once you've verified that **multipathd** is operational, run the following command:

#### multipathd -k"resize map multipath\_device"

The *multipath\_device* variable is the corresponding multipath entry of your device in /dev/ mapper. Depending on how multipathing is set up on your system, *multipath\_device* can be either of two formats:

- mpathX, where X is the corresponding entry of your device (for example, mpath0)
- a WWID; for example, 3600508b400105e210000900000490000

To determine which multipath entry corresponds to your resized logical unit, run **multipath** -11. This displays a list of all existing multipath entries in the system, along with the major and minor numbers of their corresponding devices.

Important

Do not use **multipathd** -k"resize map **multipath\_device**" if there are any commands queued to **multipath\_device**. That is, do not use this command when the **no\_path\_retry** parameter (in **/etc/multipath.conf**) is set to **"queue"**, and there are no active paths to the device.

If your system is using Red Hat Enterprise Linux 5.0-5.2, you will need to perform the following procedure to instruct the **multipathd** daemon to recognize (and adjust to) the changes you made to the resized logical unit:

Procedure 4. Resizing the Corresponding Multipath Device (Required for Red Hat Enterprise Linux 5.0 - 5.2)

1. Dump the device mapper table for the multipathed device using:

#### dmsetup table multipath\_device

2. Save the dumped device mapper table as *table\_name*. This table will be re-loaded and edited later.

3.

Examine the device mapper table. Note that the first two numbers in each line correspond to the start and end sectors of the disk, respectively.

4. Suspend the device mapper target:

#### dmsetup suspend multipath\_device

5. Open the device mapper table you saved earlier (i.e. *table\_name*). Change the second number (i.e. the disk end sector) to reflect the new number of 512 byte sectors in the disk. For example, if the new disk size is 2GB, change the second number to 4194304.

6. Reload the modified device mapper table:

#### dmsetup reload multipath\_device table\_name

7. Resume the device mapper target:

```
dmsetup resume multipath_device
```

For more information about multipathing, refer to the Using Device-Mapper Multipath<sup>9</sup> guide (in http:// www.redhat.com/docs/manuals/enterprise/).

# 15. Adding/Removing a Logical Unit Through rescan-scsibus.sh

The **sg3\_utils** package provides the **rescan-scsi-bus.sh** script, which can automatically update the logical unit configuration of the host as needed (after a device has been added to the system). The **rescan-scsi-bus.sh** script can also perform an **issue\_lip** on supported devices. For more information about how to use this script, refer to **rescan-scsi-bus.sh --help**.

To install the **sg3\_utils** package, run **yum install sg3\_utils**.

#### Known Issues With rescan-scsi-bus.sh

When using the **rescan-scsi-bus.sh** script, take note of the following known issues:

- In order for rescan-scsi-bus.sh to work properly, LUN0 must be the first mapped logical unit. The rescan-scsi-bus.sh can only detect the first mapped logical unit if it is LUN0. The rescanscsi-bus.sh will not be able to scan any other logical unit unless it detects the first mapped logical unit even if you use the --nooptscan option.
- A race condition requires that **rescan-scsi-bus.sh** be run twice if logical units are mapped for the first time. During the first scan, **rescan-scsi-bus.sh** only adds **LUNO**; all other logical units are added in the second scan.
- A bug in the **rescan-scsi-bus**. **sh** script incorrectly executes the functionality for recognizing a change in logical unit size when the **--remove** option is used.
- The rescan-scsi-bus.sh script does not recognize ISCSI logical unit removals.

## **16. Modifying Link Loss Behavior**

This section describes how to modify the link loss behavior of devices that use either fibre channel or iSCSI protocols.

## 16.1. Fibre Channel

If a driver implements the Transport **dev\_loss\_tmo** callback, access attempts to a device through a link will be blocked when a transport problem is detected. To verify if a device is blocked, run the following command:

<sup>&</sup>lt;sup>9</sup> http://www.redhat.com/docs/en-US/Red\_Hat\_Enterprise\_Linux/5.4/html/DM\_Multipath/index.html

#### cat /sys/block/device/device/state

This command will return **blocked** if the device is blocked. If the device is operating normally, this command will return **running**.

Procedure 5. Determining The State of a Remote Port

1. To determine the state of a remote port, run the following command:

#### cat /sys/class/fc\_remote\_port/rport-H:B:R/port\_state

- 2. This command will return **Blocked** when the remote port (along with devices accessed through it) are blocked. If the remote port is operating normally, the command will return **Online**.
- 3. If the problem is not resolved within **dev\_loss\_tmo** seconds, the rport and devices will be unblocked and all IO running on that device (along with any new IO sent to that device) will be failed.

#### Procedure 6. Changing dev\_loss\_tmo

To change the dev\_loss\_tmo value, echo in the desired value to the file. For example, to set dev\_loss\_tmo to 30 seconds, run:

```
echo 30 > /sys/class/fc_remote_port/rport-H:B:R/dev_loss_tmo
```

For more information about dev\_loss\_tmo, refer to Section 2.1, "Fibre Channel API".

When a device is blocked, the fibre channel class will leave the device as is; i.e. /dev/sdx will remain /dev/sdx. This is because the dev\_loss\_tmo expired. If the link problem is fixed at a later time, operations will continue using the same SCSI device and device node name.

#### Fibre Channel: remove\_on\_dev\_loss

If you prefer that devices are removed at the SCSI layer when links are marked bad (i.e. expired after dev\_loss\_tmo seconds), you can use the scsi\_transport\_fc module parameter remove\_on\_dev\_loss. When a device is removed at the SCSI layer while remove\_on\_dev\_loss is in effect, the device will be added back once all transport problems are corrected.



The use of **remove\_on\_dev\_loss** is not recommended, as removing a device at the SCSI layer does not automatically unmount any file systems from that device. When file systems from a removed device are left mounted, the device may not be properly removed from multipath or RAID devices.

Further problems may arise from this if the upper layers are not hotplug-aware. This is because the upper layers may still be holding references to the state of the device before it was originally removed. This can cause unexpected behavior when the device is added again.

#### 16.2. iSCSI Settings With dm-multipath

If dm-multipath is implemented, it is advisable to set iSCSI timers to immediately defer commands to the multipath layer. To configure this, nest the following line under device { in /etc/ multipath.conf:

features "1 queue\_if\_no\_path"

This ensures that I/O errors are retried and queued if all paths are failed in the **dm-multipath** layer.

You may need to adjust iSCSI timers further to better monitor your SAN for problems. Available iSCSI timers you can configure are *NOP-Out Interval/Timeouts* and **replacement\_timeout**, which are discussed in the following sections.

#### 16.2.1. NOP-Out Interval/Timeout

To help monitor problems the SAN, the iSCSI layer sends a NOP-Out request to each target. If a NOP-Out request times out, the iSCSI layer responds by failing any running commands and instructing the SCSI layer to requeue those commands when possible.

When **dm-multipath** is being used, the SCSI layer will fail those running commands and defer them to the multipath layer. The multipath layer then retries those commands on another path. If **dmmultipath** is *not* being used, those commands are retried five times before failing altogether.

Intervals between NOP-Out requests are 10 seconds by default. To adjust this, open **/etc/iscsi/ iscsid.conf** and edit the following line:

node.conn[0].timeo.noop\_out\_interval = [interval value]

Once set, the iSCSI layer will send a NOP-Out request to each target every [interval value] seconds.

By default, NOP-Out requests time out in 10 seconds<sup>10</sup>. To adjust this, open **/etc/iscsi/iscsid.conf** and edit the following line:

node.conn[0].timeo.noop\_out\_timeout = [timeout value]

This sets the iSCSI layer to timeout a NOP-Out request after [timeout value] seconds.

Note that this will not work for iSCSI targets already logged in or discovered once. In this case run **iscsiadm** to update configuration values explicitly.

```
iscsiadm -m node -T $target_name -p $target_ip:$port -o update -n \
node.session.timeo.replacement_timeout -v $timeout_value
```

#### **SCSI Error Handler**

 $<sup>^{\</sup>rm 10}$  Prior to Red Hat Enterprise Linux 5.4, the default NOP-Out requests time out was 15 seconds.

If the SCSI Error Handler is running, running commands on a path will not be failed immediately when a NOP-Out request times out on that path. Instead, those commands will be failed *after* **replacement\_timeout** seconds. For more information about **replacement\_timeout**, refer to *Section 16.2.2, "replacement\_timeout*".

To verify if the SCSI Error Handler is running, run:

iscsiadm -m session -P 3

#### 16.2.2. replacement\_timeout

**replacement\_timeout** controls how long the iSCSI layer should wait for a timed-out path/session to reestablish itself before failing any commands on it. The default **replacement\_timeout** value is 120 seconds.

To adjust **replacement\_timeout**, open **/etc/iscsi/iscsid.conf** and edit the following line:

node.session.timeo.replacement\_timeout = [replacement\_timeout]

The **1** queue\_if\_no\_path option in /etc/multipath.conf sets iSCSI timers to immediately defer commands to the multipath layer (refer to *Section 16.2, "iSCSI Settings With dm-multipath*"). This setting prevents I/O errors from propagating to the application; because of this, you can set **replacement\_timeout** to 15-20 seconds.

By configuring a lower **replacement\_timeout**, I/O is quickly sent to a new path and executed (in the event of a NOP-Out timeout) while the iSCSI layer attempts to re-establish the failed path/session. If all paths time out, then the multipath and device mapper layer will internally queue I/O based on the settings in **/etc/multipath.conf** instead of **/etc/iscsi/iscsid.conf**.

Note that this will not work for iSCSI targets already logged in or discovered once. In this case run **iscsiadm** to update configuration values explicitly.

iscsiadm -m node -T \$target\_name -p \$target\_ip:\$port -o update -n \
node.session.timeo.replacement\_timeout -v \$timeout\_value

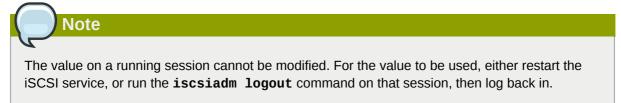


Whether your considerations are failover speed or security, the recommended value for **replacement\_timeout** will depend on other factors. These factors include the network, target, and system workload. As such, it is recommended that you thoroughly test any new configurations to **replacements\_timeout** before applying it to a mission-critical system.

The iSCSI tools do not use the **iscsid.conf** settings for portals that have already been discovered and logged into. To modify the settings of a portal that has already been discovered and set up, run:

```
iscsiadm -m node -T $target_name -p $target_ip:$port -o update -n \
node.session.timeo.replacement_timeout -v $timeout_value
```

The above command will modify the portal's record so the next time the iSCSI tools log in, that value will be used.



To set the value to be the new default for all newly discovered portals, set the value in **iscsid.conf**. Next time the **iscsiadm discovery** command is run and portals are found, the new value will be used.

#### 16.3. iSCSI Root

When accessing the root partition directly through a iSCSI disk, the iSCSI timers should be set so that iSCSI layer has several chances to try to reestablish a path/session. In addition, commands should not be quickly re-queued to the SCSI layer. This is the opposite of what should be done when **dm-multipath** is implemented.

To start with, NOP-Outs should be disabled. You can do this by setting both NOP-Out interval and timeout to zero. To set this, open **/etc/iscsi/iscsid.conf** and edit as follows:

```
node.conn[0].timeo.noop_out_interval = 0
node.conn[0].timeo.noop_out_timeout = 0
```

In line with this, **replacement\_timeout** should be set to a high number. This will instruct the system to wait a long time for a path/session to reestablish itself. To adjust **replacement\_timeout**, open / **etc/iscsi/iscsid.conf** and edit the following line:

```
node.session.timeo.replacement_timeout = replacement_timeout
```

After configuring **/etc/iscsi/iscsid.conf**, you must perform a re-discovery of the affected storage. This will allow the system to load and use any new values in **/etc/iscsi/iscsid.conf**. For more information on how to discover iSCSI devices, refer to *Section 12*, "*Scanning iSCSI Interconnects*".

#### **Configuring Timeouts for a Specific Session**

You can also configure timeouts for a specific session and make them non-persistent (instead of using **/etc/iscsi/iscsid.conf**). To do so, run the following command (replace the variables accordingly):

```
iscsiadm -m node -T target_name -p target_IP:port -o update -n
node.session.timeo.replacement_timeout -v $timeout_value
```



The configuration described here is recommended for iSCSI sessions involving root partition access. For iSCSI sessions involving access to other types of storage (namely, in systems that use dm-multipath), refer to Section 16.2, "iSCSI Settings With dm-multipath".

# **17. Controlling the SCSI Command Timer and Device Status**

The Linux SCSI layer sets a timer on each command. When this timer expires, the SCSI layer will quiesce the *host bus adapter* (HBA) and wait for all outstanding commands to either time out or complete. Afterwards, the SCSI layer will activate the driver's error handler.

When the error handler is triggered, it attempts the following operations in order (until one successfully executes):

- 1. Abort the command.
- 2. Reset the device.
- 3. Reset the bus.
- 4. Reset the host.

If all of these operations fail, the device will be set to the **offline** state. When this occurs, all IO to that device will be failed, until the problem is corrected and the user sets the device to **running**.

The process is different, however, if a device uses the fibre channel protocol and the **rport** is blocked. In such cases, the drivers wait for several seconds for the **rport** to become online again before activating the error handler. This prevents devices from becoming offline due to temporary transport problems.

#### **Device States**

To display the state of a device, use:

#### cat /sys/block/device-name/device/state

To set a device to **running** state, use:

```
echo running > /sys/block/device-name/device/state
```

#### **Command Timer**

To control the command timer, you can write to **/sys/block/***device-name*/**device/timeout**. To do so, run:

#### echo value /sys/block/device-name/device/timeout

Here, *value* is the timeout value (in seconds) you want to implement.

Alternatively, you can also modify the timeout **udev** rule. To do so, open **/etc/udev/rules.d/50-udev.rules**. You should find the following lines:

```
ACTION=="add", SUBSYSTEM=="scsi" , SYSFS{type}=="0|7|14", \
RUN+="/bin/sh -c 'echo 60 > /sys$$DEVPATH/timeout'"
```

**echo 60** refers to the timeout length, in seconds; in this case, timeout is set at 60 seconds. Replace this value with your desired timeout length.

Note that the default timeout for normal file system commands is 60 seconds when **udev** is being used. If **udev** is not in use, the default timeout is 30 seconds.

## **18. Troubleshooting**

This section provides solution to common problems users experience during online storage reconfiguration.

Logical unit removal status is not reflected on the host.

When a logical unit is deleted on a configured filer, the change is not reflected on the host. In such cases, **lvm** commands will hang indefinitely when **dm-multipath** is used, as the logical unit has now become *stale*.

To work around this, perform the following procedure:

Procedure 7. Working Around Stale Logical Units

1. Determine which **mpath** link entries in **/etc/lvm/cache/.cache** are specific to the stale logical unit. To do this, run the following command:

ls -1 /dev/mpath | grep stale-logical-unit

2. For example, if *stale-logical-unit* is 3600d0230003414f30000203a7bc41a00, the following results may appear:

lrwxrwxrwx 1 root root 7 Aug 2 10:33 /3600d0230003414f30000203a7bc41a00 -> ../dm-4 lrwxrwxrwx 1 root root 7 Aug 2 10:33 /3600d0230003414f30000203a7bc41a00p1 -> ../dm-5

This means that 3600d0230003414f30000203a7bc41a00 is mapped to two **mpath** links: **dm-4** and **dm-5**.

3. Next, open **/etc/lvm/cache/.cache**. Delete all lines containing **stale-logical-unit** and the **mpath** links that **stale-logical-unit** maps to.

Using the same example in the previous step, the lines you need to delete are:

```
/dev/dm-4
/dev/dm-5
/dev/mapper/3600d0230003414f30000203a7bc41a00
/dev/mapper/3600d0230003414f30000203a7bc41a00p1
/dev/mpath/3600d0230003414f30000203a7bc41a00
```

/dev/mpath/3600d0230003414f30000203a7bc41a00p1

## A. Revision History

**Revision 1-9** Fri Sep 16 2011 BZ#690065 Jacquelynn East jeast@redhat.com

## Index

## **Symbols**

/dev/disk persistent naming, 9

## Α

adding paths to a storage device, 13 API, fibre channel, 6 API, iSCSI, 8

## В

blocked device, verifying fibre channel modifying link loss behavior, 26

## С

changing dev\_loss\_tmo fibre channel modifying link loss behavior, 27 command timer (SCSI) Linux SCSI layer, 31 controlling SCSI command timer and device status Linux SCSI layer, 31

## D

determining remote port states fibre channel modifying link loss behavior, 27 device status Linux SCSI layer, 31 devices, removing, 11 dev\_loss\_tmo fibre channel modifying link loss behavior, 26 fibre channel API, 7 dev\_loss\_tmo, changing fibre channel modifying link loss behavior, 27 disabling NOP-Outs iSCSI configuration, 30 dm-multipath iSCSI configuration, 27 drivers (native), fibre channel, 7

#### Ε

entries, device mapper table, 25

## F

fast\_io\_fail\_tmo fibre channel API, 7 feedback contact information for this manual, 6 fibre channel API, 6 fibre channel drivers (native), 7

## Η

help getting help, 6 host fibre channel API, 7

## I

iSCSI API, 8 iSCSI logical unit, resizing, 24 iSCSI root iSCSI configuration, 30 issue\_lip fibre channel API, 7

## Μ

modifying link loss behavior, 26 fibre channel, 26

#### Ν

native fibre channel drivers, 7 NOP-Out requests modifying link loss iSCSI configuration, 28 NOP-Outs (disabling) iSCSI configuration, 30

## 0

offline status Linux SCSI layer, 31

#### Ρ

path to storage devices, adding, 13 path to storage devices, removing, 13 persistent naming, 9 port states (remote), determining fibre channel modifying link loss behavior, 27

## Q

queue\_if\_no\_path iSCSI configuration, 27 modifying link loss iSCSI configuration, 29

## R

remote port fibre channel API, 7 remote port states, determining fibre channel modifying link loss behavior, 27 remove\_on\_dev\_loss fibre channel modifying link loss behavior, 27 removing devices, 11 removing paths to a storage device, 13 replacement timeout modifying link loss iSCSI configuration, 28, 29 replacement timeoutM iSCSI configuration, 30 resized logical units, resizing, 23 resizing an iSCSI logical unit, 24 resizing multipath device resizing online resized logical units, 25 resizing resized logical units, 23 running sessions, retrieving information about **iSCSI API, 8** running status Linux SCSI layer, 31

## S

scanning storage interconnects, 15 SCSI command timer Linux SCSI layer, 31 SCSI Error Handler modifying link loss iSCSI configuration, 28 specific session timeouts, configuring iSCSI configuration, 30 storage interconnects, scanning, 15 symbolic links in /dev/disk persistent naming, 9

## Т

timeouts for a specific session, configuring iSCSI configuration, 30 transport fibre channel API, 6

#### U

udev

persistent naming, 11 udev rule (timeout) command timer (SCSI), 32 Universally Unique Identifier (UUID) persistent naming, 11 userspace API files fibre channel API, 6 UUID persistent naming, 11

## V

verifying if a device is blocked fibre channel modifying link loss behavior, 26

#### W

World Wide Identifier (WWID) persistent naming, 9 WWID persistent naming, 9