

Red Hat Enterprise Linux 5

Global File System 2

Red Hat Global File System 2



Red Hat Enterprise Linux 5 Global File System 2

Red Hat Global File System 2

Edition 7

Copyright © 2009 Red Hat Inc..

The text of and illustrations in this document are licensed by Red Hat under a Creative Commons Attribution–Share Alike 3.0 Unported license ("CC-BY-SA"). An explanation of CC-BY-SA is available at <http://creativecommons.org/licenses/by-sa/3.0/>. In accordance with CC-BY-SA, if you distribute this document or an adaptation of it, you must provide the URL for the original version.

Red Hat, as the licensor of this document, waives the right to enforce, and agrees not to assert, Section 4d of CC-BY-SA to the fullest extent permitted by applicable law.

Red Hat, Red Hat Enterprise Linux, the Shadowman logo, JBoss, MetaMatrix, Fedora, the Infinity Logo, and RHCE are trademarks of Red Hat, Inc., registered in the United States and other countries.

Linux® is the registered trademark of Linus Torvalds in the United States and other countries.

All other trademarks are the property of their respective owners.

1801 Varsity Drive
Raleigh, NC 27606-2072 USA
Phone: +1 919 754 3700
Phone: 888 733 4281
Fax: +1 919 754 3701
PO Box 13588 Research Triangle Park, NC 27709 USA

This book provides information about configuring and maintaining Red Hat GFS2 (Red Hat Global File System 2) for Red Hat Enterprise Linux 5.

Introduction	v
1. Audience	v
2. Related Documentation	v
3. Feedback	vi
4. Document Conventions	vi
4.1. Typographic Conventions	vi
4.2. Pull-quote Conventions	viii
4.3. Notes and Warnings	viii
1. GFS2 Overview	1
1.1. Before Setting Up GFS2	2
1.2. Differences between GFS and GFS2	2
1.2.1. GFS2 Command Names	2
1.2.2. Additional Differences Between GFS and GFS2	3
1.2.3. GFS2 Performance Improvements	5
2. Getting Started	7
2.1. Prerequisite Tasks	7
2.2. Initial Setup Tasks	7
3. Managing GFS2	9
3.1. Making a File System	9
3.2. Mounting a File System	12
3.3. Unmounting a File System	14
3.4. Special Considerations when Mounting GFS2 File Systems	15
3.5. GFS2 Quota Management	15
3.5.1. Setting Quotas	15
3.5.2. Displaying Quota Limits and Usage	17
3.5.3. Synchronizing Quotas	18
3.5.4. Enabling/Disabling Quota Enforcement	19
3.5.5. Enabling Quota Accounting	20
3.6. Growing a File System	20
3.7. Adding Journals to a File System	22
3.8. Data Journaling	23
3.8.1. Enabling and Disabling Data Journaling with the chattr Command	24
3.8.2. Enabling and Disabling Data Journaling with the gfs2_tool Command	24
3.9. Configuring atime Updates	25
3.9.1. Mount with relatime	26
3.9.2. Mount with noatime	26
3.10. Suspending Activity on a File System	27
3.11. Repairing a File System	27
3.12. Bind Mounts and Context-Dependent Path Names	29
3.13. Bind Mounts and File System Mount Order	30
3.14. The GFS2 Withdraw Function	32
A. Converting a File System from GFS to GFS2	33
B. Revision History	35
Index	37

Introduction

This book provides information about configuring and maintaining Red Hat GFS2 (Red Hat Global File System 2). Red Hat GFS2 can run in a single node or as part of a cluster configuration in Red Hat Cluster Suite in RHEL 5.3 and later. For information about Red Hat Cluster Suite see *Red Hat Cluster Suite Overview* and *Configuring and Managing a Red Hat Cluster*.

HTML and PDF versions of all the official Red Hat Enterprise Linux manuals and release notes are available online at <http://www.redhat.com/docs/>.

1. Audience

This book is intended primarily for Linux system administrators who are familiar with the following activities:

- Linux system administration procedures, including kernel configuration
- Installation and configuration of shared storage networks, such as Fibre Channel SANs

2. Related Documentation

For more information about using Red Hat Enterprise Linux, refer to the following resources:

- *Red Hat Enterprise Linux Installation Guide* — Provides information regarding installation of Red Hat Enterprise Linux.
- *Red Hat Enterprise Linux Deployment Guide* — Provides information regarding the deployment, configuration and administration of Red Hat Enterprise Linux 5.

For more information about Red Hat Cluster Suite, refer to the following resources:

- *Red Hat Cluster Suite Overview* — Provides a high level overview of the Red Hat Cluster Suite.
- *Configuring and Managing a Red Hat Cluster* — Provides information about installing, configuring and managing Red Hat Cluster components.
- *Logical Volume Manager Administration* — Provides a description of the Logical Volume Manager (LVM), including information on running LVM in a clustered environment.
- *Global File System: Configuration and Administration* — Provides information about installing, configuring, and maintaining Red Hat GFS (Red Hat Global File System).
- *Using Device-Mapper Multipath* — Provides information about using the Device-Mapper Multipath feature of Red Hat Enterprise Linux.
- *Using GNBD with Global File System* — Provides an overview on using Global Network Block Device (GNBD) with Red Hat GFS.
- *Linux Virtual Server Administration* — Provides information on configuring high-performance systems and services with the Linux Virtual Server (LVS).
- *Red Hat Cluster Suite Release Notes* — Provides information about the current release of Red Hat Cluster Suite.

Red Hat Cluster Suite documentation and other Red Hat documents are available in HTML, PDF, and RPM versions on the Red Hat Enterprise Linux Documentation CD and online at <http://www.redhat.com/docs/>.

3. Feedback

If you spot a typo, 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/bugzilla/>) against the component **rh-cs**.

Be sure to mention the manual's identifier:

```
rh-gfs2(EN)-5 (2010-02-03T15:15)
```

By mentioning this manual's identifier, we know exactly which version of the guide you have.

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

4. 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](#)¹ 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.

4.1. 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:

¹ <https://fedorahosted.org/liberation-fonts/>

Press **Enter** to execute the command.

Press **Ctrl+Alt+F1** to switch to the first virtual terminal. Press **Ctrl+Alt+F7** 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).

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** → **Preferences** → **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** → **Accessories** → **Character Map** from the main menu bar. Next, choose **Search** → **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 **Edit** → **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.

4.2. 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:

```
books      Desktop  documentation  drafts  mss    photos  stuff  svn
books_tests Desktop1  downloads      images  notes  scripts svgs
```

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"));
    }
}
```

4.3. 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' won't cause data loss but may cause irritation and frustration.



Warning

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

GFS2 Overview

The Red Hat GFS2 file system is a native file system that interfaces directly with the Linux kernel file system interface (VFS layer). A GFS2 file system can be implemented in a standalone system or as part of a cluster configuration. When implemented as a cluster file system, GFS2 employs distributed metadata and multiple journals.



Note

Red Hat does not support using GFS2 for cluster file system deployments greater than 16 nodes.

GFS2 is based on a 64-bit architecture, which can theoretically accommodate an 8 EB file system. However, the current supported maximum size of a GFS2 file system is 25 TB. If your system requires GFS2 file systems larger than 25 TB, contact your Red Hat service representative.

When determining the size of your file system, you should consider your recovery needs. Running the **fsck** command on a very large file system can take a long time and consume a large amount of memory. Additionally, in the event of a disk or disk-subsystem failure, recovery time is limited by the speed of your backup media.

When configured in a Red Hat Cluster Suite, Red Hat GFS2 nodes can be configured and managed with Red Hat Cluster Suite configuration and management tools. Red Hat GFS2 then provides data sharing among GFS2 nodes in a Red Hat cluster, with a single, consistent view of the file system name space across the GFS2 nodes. This allows processes on different nodes to share GFS2 files in the same way that processes on the same node can share files on a local file system, with no discernible difference. For information about Red Hat Cluster Suite refer to *Configuring and Managing a Red Hat Cluster*.

A GFS2 file system must be created on an LVM logical volume that is a linear or mirrored volume. LVM logical volumes in a Red Hat Cluster suite are managed with CLVM, which is a cluster-wide implementation of LVM, enabled by the CLVM daemon, **clvmd**, running in a Red Hat Cluster Suite cluster. The daemon makes it possible to use LVM2 to manage logical volumes across a cluster, allowing all nodes in the cluster to share the logical volumes. For information on the LVM volume manager, see *Logical Volume Manager Administration*.

The **gfs2.ko** kernel module implements the GFS2 file system and is loaded on GFS2 cluster nodes.



Note

When you configure a GFS2 file system as a cluster file system, you must ensure that all nodes in the cluster have access to the the shared file system. Asymmetric cluster configurations in which some nodes have access to the file system and others do not are not supported.

This chapter provides some basic, abbreviated information as background to help you understand GFS2. It contains the following sections:

- [Section 1.1, “Before Setting Up GFS2”](#)
- [Section 1.2, “Differences between GFS and GFS2”](#)

1.1. Before Setting Up GFS2

Before you install and set up GFS2, note the following key characteristics of your GFS2 file systems:

GFS2 nodes

Determine which nodes in the Red Hat Cluster Suite will mount the GFS2 file systems.

Number of file systems

Determine how many GFS2 file systems to create initially. (More file systems can be added later.)

File system name

Determine a unique name for each file system. The name must be unique for all **lock_dlm** filesystems over the cluster, and for all filesystems (**lock_dlm** and **lock_nolock**) on each local node. Each file system name is required in the form of a parameter variable. For example, this book uses file system names **mydata1** and **mydata2** in some example procedures.

Journals

Determine the number of journals for your GFS2 file systems. One journal is required for each node that mounts a GFS2 file system. GFS2 allows you to add journals dynamically at a later point as additional servers mount a filesystem. For information on adding journals to a GFS2 file system, see [Section 3.7, “Adding Journals to a File System”](#).

GNBD server nodes

If you are using GNBD, determine how many GNBD server nodes are needed. Note the hostname and IP address of each GNBD server node for setting up GNBD clients later. For information on using GNBD with GFS2, see the *Using GNBD with Global File System* document.

Storage devices and partitions

Determine the storage devices and partitions to be used for creating logical volumes (via CLVM) in the file systems.



Note

You may see performance problems with GFS2 when many create and delete operations are issued from more than one node in the same directory at the same time. If this causes performance problems in your system, you should localize file creation and deletions by a node to directories specific to that node as much as possible.

1.2. Differences between GFS and GFS2

This section lists the improvements and changes that GFS2 offers over GFS.

Migrating from GFS to GFS2 requires that you convert your GFS file systems to GFS2 with the **gfs2_convert** utility. For information on the **gfs2_convert** utility, see [Appendix A, Converting a File System from GFS to GFS2](#).

1.2.1. GFS2 Command Names

In general, the functionality of GFS2 is identical to GFS. The names of the file system commands, however, specify GFS2 instead of GFS. [Table 1.1, “GFS and GFS2 Commands”](#) shows the equivalent GFS and GFS2 commands.

GFS Command	GFS2 Command	Description
mount	mount	Mount a file system. The system can determine whether the file system is a GFS or GFS2 file system type. For information on the GFS2 mount options see the <code>gfs2_mount(8)</code> man page.
umount	umount	Unmount a file system.
fsck gfs_fsck	fsck fsck.gfs2	Check and repair an unmounted file system.
gfs_grow	gfs2_grow	Grow a mounted file system.
gfs_jadd	gfs2_jadd	Add a journal to a mounted file system
gfs_mkfs mkfs -t gfs	mkfs.gfs2 mkfs -t gfs2	Create a file system on a storage device.
gfs_quota	gfs2_quota	Manage quotas on a mounted file system.
gfs_tool	gfs2_tool	Configure, tune, or gather information about a file system.
gfs_edit	gfs2_edit	Display, print, or edit file system internal structures. The gfs2_edit command can be used for GFS file systems as well as GFS2 file system.

Table 1.1. GFS and GFS2 Commands

For a full listing of the supported options for the GFS2 file system commands, see the man pages for those commands.

1.2.2. Additional Differences Between GFS and GFS2

This section summarizes the additional differences in GFS and GFS2 administration that are not described in [Section 1.2.1, “GFS2 Command Names”](#).

Context-Dependent Path Names

GFS2 file systems do not provide support for context-dependent path names, which allow you to create symbolic links that point to variable destination files or directories. For this functionality in GFS2, you can use the **bind** option of the **mount** command. For information on managing pathnames in GFS2, see [Section 3.12, “Bind Mounts and Context-Dependent Path Names”](#).

gfs2.ko Module

The kernel module that implements the GFS file system is **gfs.ko**. The kernel module that implements the GFS2 file system is **gfs2.ko**.

Enabling Quota Enforcement in GFS2

In GFS2 file systems, quota enforcement is disabled by default and must be explicitly enabled. To enable and disable quotas for GFS2 file systems, you use the **quota=on|off|account** option for the **mount** command. For information on enabling and disabling quota enforcement, see [Section 3.5.4, “Enabling/Disabling Quota Enforcement”](#).

Data Journaling

GFS2 file systems support the use of the **chattr** command to set and clear the **j** flag on a file or directory. Setting the **+j** flag on a file enables data journaling on that file. Setting the **+j** flag on a

directory means "inherit jdata", which indicates that all files and directories subsequently created in that directory are journaled. Using the **chattr** command is the preferred way to enable and disable data journaling on a file.

Adding Journals Dynamically

In GFS file systems, journals are embedded metadata that exists outside of the file system, making it necessary to extend the size of the logical volume that contains the file system before adding journals. In GFS2 file systems, journals are plain (though hidden) files. This means that for GFS2 file systems, journals can be dynamically added as additional servers mount a filesystem, as long as space remains on the file system for the additional journals. For information on adding journals to a GFS2 file system, see [Section 3.7, "Adding Journals to a File System"](#).

atime_quantum parameter removed

The GFS2 file system does not support the **atime_quantum** tunable parameter, which can be used by the GFS file system to specify how often **atime** updates occur. In its place GFS2 supports the **relatime** and **noatime** mount options. The **relatime** mount option is recommended to achieve similar behavior to setting the **atime_quantum** parameter in GFS.

The data= option of the mount command

When mounting GFS2 file systems, you can specify the **data=ordered** or **data=writeback** option of the **mount**. When **data=ordered** is set, the user data modified by a transaction is flushed to the disk before the transaction is committed to disk. This should prevent the user from seeing uninitialized blocks in a file after a crash. When **data=writeback** is set, the user data is written to the disk at any time after it is dirtied. This does not provide the same consistency guarantee as **ordered** mode, but it should be slightly faster for some workloads. The default is **ordered** mode.

The gfs2_tool command

The **gfs2_tool** command supports a different set of options for GFS2 than the **gfs_tool** command supports for GFS:

- The **gfs2_tool** command supports a **journals** parameter that prints out information about the currently configured journal, including how many journals the file system contains.
- The **gfs2_tool** command does not support the **counters** flag, which the **gfs_tool** command uses to display GFS statistics.
- The **gfs2_tool** command does not support the **inherit_jdata** flag. To flag a directory as "inherit jdata", you can set the **jdata** flag on the directory or you can use the **chattr** command to set the **+j** flag on the directory. Using the **chattr** command is the preferred way to enable and disable data journaling on a file.

The gfs2_edit command

The **gfs2_edit** command supports a different set of options for GFS2 than the **gfs_edit** command supports for GFS.

1.2.3. GFS2 Performance Improvements

There are many features of GFS2 filesystems that do not result in a difference in the user interface from GFS file systems but which improve file system performance.

A GFS2 filesystem provides improved file system performance in the following ways:

- Better performance for heavy usage in a single directory.
- Faster synchronous I/O operations
- Faster cached reads (no locking overhead)
- Faster direct I/O with preallocated files (provided I/O size is reasonably large, such as 4M blocks)
- Faster I/O operations in general
- Execution of the **df** command is much faster, because of faster **statfs** calls.
- The **atime** mode has been improved to reduce the number of write I/O operations generated by **atime** when compared with GFS.

GFS2 file system provide broader and more mainstream support in the following ways.

- GFS2 is part of the upstream kernel (integrated into 2.6.19).
- GFS2 supports the following features:
 - SELinux extended attributes.
 - the **lsattr()** and **chattr()** attribute settings via standard **ioctl()** calls.
 - nanosecond timestamps

A GFS2 file system provides the following improvements to the internal efficiency of the file system.

- GFS2 uses less kernel memory
- GFS2 requires no metadata generation numbers.

Allocating GFS2 metadata does not require reads. Copies of metadata blocks in multiple journals are managed by revoking blocks from the journal before lock release.

- GFS2 includes a much simpler log manager that knows nothing about unlinked inodes or quota changes.
- The **gfs2_grow** and **gfs2_jadd** commands use locking to prevent multiple instances running at the same time.
- The ACL code on has been simplified for calls like **creat()** and **mkdir()**.
- Unlinked inodes, quota changes, and **statfs** changes are recovered without remounting the journal.

Getting Started

This chapter describes procedures for initial setup of GFS2 and contains the following sections:

- [Section 2.1, “Prerequisite Tasks”](#)
- [Section 2.2, “Initial Setup Tasks”](#)

2.1. Prerequisite Tasks

Before setting up Red Hat GFS2, make sure that you have noted the key characteristics of the GFS2 nodes (refer to [Section 1.1, “Before Setting Up GFS2”](#)). Also, make sure that the clocks on the GFS2 nodes are synchronized. It is recommended that you use the Network Time Protocol (NTP) software provided with your Red Hat Enterprise Linux distribution.



Note

The system clocks in GFS2 nodes must be within a few minutes of each other to prevent unnecessary inode time-stamp updating. Unnecessary inode time-stamp updating severely impacts cluster performance.

2.2. Initial Setup Tasks

Initial GFS2 setup consists of the following tasks:

1. Setting up logical volumes.
2. Making a GFS2 file system.
3. Mounting file systems.

Follow these steps to set up GFS2 initially.

1. Using LVM, create a logical volume for each Red Hat GFS2 file system.



Note

You can use `init.d` scripts included with Red Hat Cluster Suite to automate activating and deactivating logical volumes. For more information about `init.d` scripts, refer to *Configuring and Managing a Red Hat Cluster*.

2. Create GFS2 file systems on logical volumes created in Step 1. Choose a unique name for each file system. For more information about creating a GFS2 file system, refer to [Section 3.1, “Making a File System”](#).

You can use either of the following formats to create a clustered GFS2 file system:

```
mkfs.gfs2 -p lock_dlm -t ClusterName:FSName -j NumberJournals BlockDevice
```

```
mkfs -t gfs2 -p lock_dlm -t LockTableName -j NumberJournals BlockDevice
```

You can use either of the following formats to create a local GFS2 file system:

```
mkfs.gfs2 -p lock_nolock -j NumberJournals BlockDevice
```

```
mkfs -t gfs2 -p lock_nolock -j NumberJournals BlockDevice
```

For more information on creating a GFS2 file system, see [Section 3.1, “Making a File System”](#).

3. At each node, mount the GFS2 file systems. For more information about mounting a GFS2 file system, see [Section 3.2, “Mounting a File System”](#).

Command usage:

```
mount BlockDevice MountPoint
```

```
mount -o acl BlockDevice MountPoint
```

The **-o acl** mount option allows manipulating file ACLs. If a file system is mounted without the **-o acl** mount option, users are allowed to view ACLs (with **getfacl**), but are not allowed to set them (with **setfacl**).



Note

You can use **init.d** scripts included with Red Hat Cluster Suite to automate mounting and unmounting GFS2 file systems. For more information about **init.d** scripts, refer to *Configuring and Managing a Red Hat Cluster*.

Managing GFS2

This chapter describes the tasks and commands for managing GFS2 and consists of the following sections:

- [Section 3.1, “Making a File System”](#)
- [Section 3.2, “Mounting a File System”](#)
- [Section 3.3, “Unmounting a File System”](#)
- [Section 3.5, “GFS2 Quota Management”](#)
- [Section 3.6, “Growing a File System”](#)
- [Section 3.7, “Adding Journals to a File System”](#)
- [Section 3.8, “Data Journaling”](#)
- [Section 3.9, “Configuring *atime* Updates”](#)
- [Section 3.10, “Suspending Activity on a File System”](#)
- [Section 3.11, “Repairing a File System”](#)
- [Section 3.12, “Bind Mounts and Context-Dependent Path Names”](#)
- [Section 3.13, “Bind Mounts and File System Mount Order”](#)
- [Section 3.14, “The GFS2 Withdraw Function”](#)

3.1. Making a File System

You create a GFS2 file system with the **mkfs.gfs2** command. You can also use the **mkfs** command with the **-t gfs2** option specified. A file system is created on an activated LVM volume. The following information is required to run the **mkfs.gfs2** command:

- Lock protocol/module name (the lock protocol for a cluster is **lock_dlm**)
- Cluster name (when running as part of a cluster configuration)
- Number of journals (one journal required for each node that may be mounting the file system)

When creating a GFS2 file system, you can use the **mkfs.gfs2** command directly, or you can use the **mkfs** command with the **-t** parameter specifying a filesystem of type **gfs2**, followed by the gfs2 file system options.



Note

Once you have created a GFS2 file system with the **mkfs.gfs2** command, you cannot decrease the size of the file system. You can, however, increase the size of an existing file system with the **gfs2_grow** command, as described in [Section 3.6, “Growing a File System”](#).

Usage

When creating a clustered GFS2 filesystem, you can use either of the following formats:

```
mkfs.gfs2 -p LockProtoName -t LockTableName -j NumberJournals BlockDevice
```

```
mkfs -t gfs2 -p LockProtoName -t LockTableName -j NumberJournals BlockDevice
```

When creating a local GFS2 filesystem, you can use either of the following formats:

```
mkfs.gfs2 -p LockProtoName -j NumberJournals BlockDevice
```

```
mkfs -t gfs2 -p LockProtoName -j NumberJournals BlockDevice
```



Warning

Make sure that you are very familiar with using the *LockProtoName* and *LockTableName* parameters. Improper use of the *LockProtoName* and *LockTableName* parameters may cause file system or lock space corruption.

LockProtoName

Specifies the name of the locking protocol to use. The lock protocol for a cluster is **lock_dlm**.

LockTableName

This parameter is specified for GFS2 filesystem in a cluster configuration. It has two parts separated by a colon (no spaces) as follows: *ClusterName:FSName*

- *ClusterName*, the name of the Red Hat cluster for which the GFS2 file system is being created.
- *FSName*, the file system name, can be 1 to 16 characters long. The name must be unique for all **lock_dlm** filesystems over the cluster, and for all filesystems (**lock_dlm** and **lock_nolock**) on each local node.

Number

Specifies the number of journals to be created by the **mkfs.gfs2** command. One journal is required for each node that mounts the file system. For GFS2 file systems, more journals can be added later without growing the filesystem, as described in [Section 3.7, “Adding Journals to a File System”](#).

BlockDevice

Specifies a logical or physical volume.

Examples

In these example, **lock_dlm** is the locking protocol that the file system uses, since this is a clustered file system. The cluster name is **alpha**, and the file system name is **mydata1**. The file system contains eight journals and is created on **/dev/vg01/lvol0**.

```
mkfs.gfs2 -p lock_dlm -t alpha:mydata1 -j 8 /dev/vg01/lvol0
```

```
mkfs -t gfs2 -p lock_dlm -t alpha:mydata1 -j 8 /dev/vg01/lvol0
```

In these examples, a second **lock_dlm** file system is made, which can be used in cluster **alpha**. The file system name is **mydata2**. The file system contains eight journals and is created on **/dev/vg01/lvol1**.

```
mkfs.gfs2 -p lock_dlm -t alpha:mydata2 -j 8 /dev/vg01/lvol1
```

```
mkfs -t gfs2 -p lock_dlm -t alpha:mydata2 -j 8 /dev/vg01/lvol1
```

Complete Options

Table 3.1, “Command Options: **mkfs.gfs2**” describes the **mkfs.gfs2** command options (flags and parameters).

Flag	Parameter	Description
-c	<i>Megabytes</i>	Sets the initial size of each journal's quota change file to <i>Megabytes</i> .
-D		Enables debugging output.
-h		Help. Displays available options.
-J	<i>MegaBytes</i>	Specifies the size of the journal in megabytes. Default journal size is 128 megabytes. The minimum size is 8 megabytes. Larger journals improve performance, although they use more memory than smaller journals.
-j	<i>Number</i>	Specifies the number of journals to be created by the mkfs.gfs2 command. One journal is required for each node that mounts the file system. If this option is not specified, one journal will be created. For GFS2 file systems, you can add additional journals at a later time without growing the file system.
-O		Prevents the mkfs.gfs2 command from asking for confirmation before writing the file system.
-p	<i>LockProtoName</i>	Specifies the name of the locking protocol to use. Recognized locking protocols include: lock_dlm — The standard locking module, required for a clustered file system. lock_nolock — Used when GFS2 is acting as a local file system (one node only).
-q		Quiet. Do not display anything.
-r	<i>MegaBytes</i>	Specifies the size of the resource groups in megabytes. The minimum resource group size is 32 MB. The maximum resource group size is 2048 MB. A large

Flag	Parameter	Description
		resource group size size may increase performance on very large file systems. If this is not specified, <code>mkfs.gfs2</code> chooses the resource group size based on the size of the file system: average size file systems will have 256 MB resource groups, and bigger file systems will have bigger RGs for better performance.
-t	<i>LockTableName</i>	A unique identifier that specifies the lock table field when you use the lock_dlm protocol; the lock_nolock protocol does not use this parameter. This parameter has two parts separated by a colon (no spaces) as follows: <i>ClusterName:FSName</i> . <i>ClusterName</i> is the name of the Red Hat cluster for which the GFS2 file system is being created; only members of this cluster are permitted to use this file system. The cluster name is set in the <code>/etc/cluster/cluster.conf</code> file via the Cluster Configuration Tool and displayed at the Cluster Status Tool in the Red Hat Cluster Suite cluster management GUI. <i>FSName</i> , the file system name, can be 1 to 16 characters in length, and the name must be unique among all file systems in the cluster.
-u	<i>MegaBytes</i>	Specifies the initial size of each journal's unlinked tag file.
-v		Displays command version information.

Table 3.1. Command Options: `mkfs.gfs2`

3.2. Mounting a File System

Before you can mount a GFS2 file system, the file system must exist (refer to [Section 3.1, “Making a File System”](#)), the volume where the file system exists must be activated, and the supporting clustering and locking systems must be started (refer to *Configuring and Managing a Red Hat Cluster*). After those requirements have been met, you can mount the GFS2 file system as you would any Linux file system.

To manipulate file ACLs, you must mount the file system with the **-o acl** mount option. If a file system is mounted without the **-o acl** mount option, users are allowed to view ACLs (with **getfacl**), but are not allowed to set them (with **setfacl**).

Usage

Mounting Without ACL Manipulation

```
mount BlockDevice MountPoint
```

Mounting With ACL Manipulation

```
mount -o acl BlockDevice MountPoint
```

-o acl

GFS2-specific option to allow manipulating file ACLs.

BlockDevice

Specifies the block device where the GFS2 file system resides.

MountPoint

Specifies the directory where the GFS2 file system should be mounted.

Example

In this example, the GFS2 file system on **/dev/vg01/lvol0** is mounted on the **/mygfs2** directory.

```
mount /dev/vg01/lvol0 /mygfs2
```

Complete Usage

```
mount BlockDevice MountPoint -o option
```

The **-o option** argument consists of GFS2-specific options (refer to [Table 3.2, “GFS2-Specific Mount Options”](#)) or acceptable standard Linux **mount -o** options, or a combination of both. Multiple *option* parameters are separated by a comma and no spaces.



Note

The **mount** command is a Linux system command. In addition to using GFS2-specific options described in this section, you can use other, standard, **mount** command options (for example, **-r**). For information about other Linux **mount** command options, see the Linux **mount** man page.

[Table 3.2, “GFS2-Specific Mount Options”](#) describes the available GFS2-specific **-o option** values that can be passed to GFS2 at mount time.

Option	Description
acl	Allows manipulating file ACLs. If a file system is mounted without the acl mount option, users are allowed to view ACLs (with getfacl), but are not allowed to set them (with setfacl).
data=[ordered writeback]	When data=ordered is set, the user data modified by a transaction is flushed to the disk before the transaction is committed to disk. This should prevent the user from seeing uninitialized blocks in a file after a crash. When data=writeback mode is set, the user data is written to the disk at any time after it is dirtied; this does not provide the same consistency guarantee

Option	Description
	as ordered mode, but it should be slightly faster for some workloads. The default value is ordered mode.
ignore_local_fs Caution: This option should <i>not</i> be used when GFS2 file systems are shared.	Forces GFS2 to treat the file system as a multihost file system. By default, using lock_nolock automatically turns on the localcaching and localflocks flags.
localcaching Caution: This option should <i>not</i> be used when GFS2 file systems are shared.	Tells GFS2 that it is running as a local file system. GFS2 can then turn on selected optimization capabilities that are not available when running in cluster mode. The localcaching flag is automatically turned on by lock_nolock .
localflocks Caution: This option should not be used when GFS2 file systems are shared.	Tells GFS2 to let the VFS (virtual file system) layer do all flock and fcntl. The localflocks flag is automatically turned on by lock_nolock .
lockproto=LockModuleName	Allows the user to specify which locking protocol to use with the file system. If <i>LockModuleName</i> is not specified, the locking protocol name is read from the file system superblock.
locktable=LockTableName	Allows the user to specify which locking table to use with the file system.
quota=[off/account/on]	Turns quotas on or off for a file system. Setting the quotas to be in the account state causes the per UID/GID usage statistics to be correctly maintained by the file system; limit and warn values are ignored. The default value is off .
upgrade	Upgrade the on-disk format of the file system so that it can be used by newer versions of GFS2.

Table 3.2. GFS2-Specific Mount Options

3.3. Unmounting a File System

The GFS2 file system can be unmounted the same way as any Linux file system — by using the **umount** command.



Note

The **umount** command is a Linux system command. Information about this command can be found in the Linux **umount** command man pages.

Usage

```
umount MountPoint
```

MountPoint

Specifies the directory where the GFS2 file system is currently mounted.

3.4. Special Considerations when Mounting GFS2 File Systems

GFS2 file systems that have been mounted manually rather than automatically through an entry in the **fstab** file will not be known to the system when file systems are unmounted at system shutdown. As a result, the GFS2 script will not unmount the GFS2 file system. After the GFS2 shutdown script is run, the standard shutdown process kills off all remaining user processes, including the cluster infrastructure, and tries to unmount the filesystem. This unmount will fail without the cluster infrastructure and the system will hang.

To prevent the system from hanging when the GFS2 file systems are unmounted, you should do one of the following:

- Always use an entry in the **fstab** file to mount the GFS2 file system.
- If a GFS2 file system has been mounted manually with the **mount** command, be sure to unmount the file system manually with the **umount** command before rebooting or shutting down the system.

If your file system hangs while it is being unmounted during system shutdown under these circumstances, perform a hardware reboot. It is unlikely that any data will be lost since the file system is synced earlier in the shutdown process.

3.5. GFS2 Quota Management

File-system quotas are used to limit the amount of file system space a user or group can use. A user or group does not have a quota limit until one is set. GFS2 keeps track of the space used by each user and group even when there are no limits in place. GFS2 updates quota information in a transactional way so system crashes do not require quota usages to be reconstructed.

To prevent a performance slowdown, a GFS2 node synchronizes updates to the quota file only periodically. The "fuzzy" quota accounting can allow users or groups to slightly exceed the set limit. To minimize this, GFS2 dynamically reduces the synchronization period as a "hard" quota limit is approached.

GFS2 uses its **gfs2_quota** command to manage quotas. Other Linux quota facilities cannot be used with GFS2.

3.5.1. Setting Quotas

Two quota settings are available for each user ID (UID) or group ID (GID): a *hard limit* and a *warn limit*.

A hard limit is the amount of space that can be used. The file system will not let the user or group use more than that amount of disk space. A hard limit value of *zero* means that no limit is enforced.

A warn limit is usually a value less than the hard limit. The file system will notify the user or group when the warn limit is reached to warn them of the amount of space they are using. A warn limit value of *zero* means that no limit is enforced.

Limits are set using the **gfs2_quota** command. The command only needs to be run on a single node where GFS2 is mounted.

By default, quota enforcement is not set on GFS2 file systems. To enable quota accounting, use the **quota=** of the **mount** command when mounting the GFS2 file system, as described in [Section 3.5.4, “Enabling/Disabling Quota Enforcement”](#).

Usage

Setting Quotas, Hard Limit

```
gfs2_quota limit -u User -l Size -f MountPoint
```

```
gfs2_quota limit -g Group -l Size -f MountPoint
```

Setting Quotas, Warn Limit

```
gfs2_quota warn -u User -l Size -f MountPoint
```

```
gfs2_quota warn -g Group -l Size -f MountPoint
```

User

A user ID to limit or warn. It can be either a user name from the password file or the UID number.

Group

A group ID to limit or warn. It can be either a group name from the group file or the GID number.

Size

Specifies the new value to limit or warn. By default, the value is in units of megabytes. The additional **-k**, **-s** and **-b** flags change the units to kilobytes, sectors, and file system blocks, respectively.

MountPoint

Specifies the GFS2 file system to which the actions apply.

Examples

This example sets the hard limit for user *Bert* to 1024 megabytes (1 gigabyte) on file system **/mygfs2**.

```
gfs2_quota limit -u Bert -l 1024 -f /mygfs2
```

This example sets the warn limit for group ID 21 to 50 kilobytes on file system **/mygfs2**.

```
gfs2_quota warn -g 21 -l 50 -k -f /mygfs2
```

3.5.2. Displaying Quota Limits and Usage

Quota limits and current usage can be displayed for a specific user or group using the **gfs2_quota get** command. The entire contents of the quota file can also be displayed using the **gfs2_quota list** command, in which case all IDs with a non-zero hard limit, warn limit, or value are listed.

Usage

Displaying Quota Limits for a User

```
gfs2_quota get -u User -f MountPoint
```

Displaying Quota Limits for a Group

```
gfs2_quota get -g Group -f MountPoint
```

Displaying Entire Quota File

```
gfs2_quota list -f MountPoint
```

User

A user ID to display information about a specific user. It can be either a user name from the password file or the UID number.

Group

A group ID to display information about a specific group. It can be either a group name from the group file or the GID number.

MountPoint

Specifies the GFS2 file system to which the actions apply.

Command Output

GFS2 quota information from the **gfs2_quota** command is displayed as follows:

```
user User: limit:LimitSize warn:WarnSize value:Value
group Group: limit:LimitSize warn:WarnSize value:Value
```

The *LimitSize*, *WarnSize*, and *Value* numbers (values) are in units of megabytes by default. Adding the **-k**, **-s**, or **-b** flags to the command line change the units to kilobytes, sectors, or file system blocks, respectively.

User

A user name or ID to which the data is associated.

Group

A group name or ID to which the data is associated.

LimitSize

The hard limit set for the user or group. This value is zero if no limit has been set.

Value

The actual amount of disk space used by the user or group.

Comments

When displaying quota information, the **gfs2_quota** command does not resolve UIDs and GIDs into names if the **-n** option is added to the command line.

Space allocated to GFS2's hidden files can be left out of displayed values for the root UID and GID by adding the **-d** option to the command line. This is useful when trying to match the numbers from **gfs2_quota** with the results of a **du** command.

Examples

This example displays quota information for all users and groups that have a limit set or are using any disk space on file system **/mygfs2**.

```
gfs2_quota list -f /mygfs2
```

This example displays quota information in sectors for group **users** on file system **/mygfs2**.

```
gfs2_quota get -g users -f /mygfs2 -s
```

3.5.3. Synchronizing Quotas

GFS2 stores all quota information in its own internal file on disk. A GFS2 node does not update this quota file for every file system write; rather, it updates the quota file once every 60 seconds. This is necessary to avoid contention among nodes writing to the quota file, which would cause a slowdown in performance.

As a user or group approaches their quota limit, GFS2 dynamically reduces the time between its quota-file updates to prevent the limit from being exceeded. The normal time period between quota synchronizations is a tunable parameter, **quota_quantum**, and can be changed using the **gfs2_tool** command. By default, the time period is 60 seconds. Also, the **quota_quantum** parameter must be set on each node and each time the file system is mounted. (Changes to the **quota_quantum** parameter are not persistent across unmounts.)

You can use the **gfs2_quota sync** command to synchronize the quota information from a node to the on-disk quota file between the automatic updates performed by GFS2.

Usage

Synchronizing Quota Information

```
gfs2_quota sync -f MountPoint
```

MountPoint

Specifies the GFS2 file system to which the actions apply.

Tuning the Time Between Synchronizations

```
gfs2_tool settune MountPoint quota_quantum Seconds
```

MountPoint

Specifies the GFS2 file system to which the actions apply.

Seconds

Specifies the new time period between regular quota-file synchronizations by GFS2. Smaller values may increase contention and slow down performance.

Examples

This example synchronizes the quota information from the node it is run on to file system **/mygfs2**.

```
gfs2_quota sync -f /mygfs2
```

This example changes the default time period between regular quota-file updates to one hour (3600 seconds) for file system **/mygfs2** on a single node.

```
gfs2_tool settune /mygfs2 quota_quantum 3600
```

3.5.4. Enabling/Disabling Quota Enforcement

In GFS2 file systems, quota enforcement is disabled by default. To enable quota enforcement for a file system, mount the file system with the **quota=on** option specified.

Usage

```
mount -o quota=on BlockDevice MountPoint
```

To mount a file system with quota enforcement disabled, mount the file system with the **quota=off** option specified. This is the default setting.

```
mount -o quota=off BlockDevice MountPoint
```

-o quota={on|off}

Specifies that quota enforcement is enabled or disabled when the file system is mounted.

BlockDevice

Specifies the block device where the GFS2 file system resides.

MountPoint

Specifies the directory where the GFS2 file system should be mounted.

Examples

In this example, the GFS2 file system on **/dev/vg01/lvo10** is mounted on the **/mygfs2** directory with quota enforcement enabled.

```
mount -o quota=on /dev/vg01/lvol0 /mygfs2
```

3.5.5. Enabling Quota Accounting

It is possible to keep track of disk usage and maintain quota accounting for every user and group without enforcing the limit and warn values. To do this, mount the file system with the **quota=account** option specified.

Usage

```
mount -o quota=account BlockDevice MountPoint
```

-o quota=account

Specifies that user and group usage statistics are maintained by the file system, even though the quota limits are not enforced.

BlockDevice

Specifies the block device where the GFS2 file system resides.

MountPoint

Specifies the directory where the GFS2 file system should be mounted.

Example

In this example, the GFS2 file system on **/dev/vg01/lvol0** is mounted on the **/mygfs2** directory with quota accounting enabled.

```
mount -o quota=account /dev/vg01/lvol0 /mygfs2
```

3.6. Growing a File System

The **gfs2_grow** command is used to expand a GFS2 file system after the device where the file system resides has been expanded. Running a **gfs2_grow** command on an existing GFS2 file system fills all spare space between the current end of the file system and the end of the device with a newly initialized GFS2 file system extension. When the fill operation is completed, the resource index for the file system is updated. All nodes in the cluster can then use the extra storage space that has been added.

The **gfs2_grow** command must be run on a mounted file system, but only needs to be run on one node in a cluster. All the other nodes sense that the expansion has occurred and automatically start using the new space.



Note

Once you have created a GFS2 file system with the **mkfs.gfs2** command, you cannot decrease the size of the file system.

Usage

```
gfs2_grow MountPoint
```

MountPoint

Specifies the GFS2 file system to which the actions apply.

Comments

Before running the **gfs2_grow** command:

- Back up important data on the file system.
- Determine the volume that is used by the file system to be expanded by running a **df *MountPoint*** command.
- Expand the underlying cluster volume with LVM. For information on administering LVM volumes, see the *LVM Administrator's Guide*

After running the **gfs2_grow** command, run a **df** command to check that the new space is now available in the file system.

Examples

In this example, the file system on the **/mygfs2fs** directory is expanded.

```
[root@dash-01 ~]# gfs2_grow /mygfs2fs
FS: Mount Point: /mygfs2fs
FS: Device:      /dev/mapper/gfs2testvg-gfs2testlv
FS: Size:        524288 (0x80000)
FS: RG size:     65533 (0xffffd)
DEV: Size:       655360 (0xa0000)
The file system grew by 512MB.
gfs2_grow complete.
```

Complete Usage

```
gfs2_grow [Options] {MountPoint | Device} [MountPoint | Device]
```

MountPoint

Specifies the directory where the GFS2 file system is mounted.

Device

Specifies the device node of the file system.

[Table 3.3, “GFS2-specific Options Available While Expanding A File System”](#) describes the GFS2-specific options that can be used while expanding a GFS2 file system.

Option	Description
-h	Help. Displays a short usage message.

Option	Description
-q	Quiet. Turns down the verbosity level.
-r MegaBytes	Specifies the size of the new resource group. The default size is 256MB.
-T	Test. Do all calculations, but do not write any data to the disk and do not expand the file system.
-V	Displays command version information.

Table 3.3. GFS2-specific Options Available While Expanding A File System

3.7. Adding Journals to a File System

The **gfs2_jadd** command is used to add journals to a GFS2 file system. You can add journals to a GFS2 file system dynamically at any point without expanding the underlying logical volume. The **gfs2_jadd** command must be run on mounted a file system, but it needs to be run on only one node in the cluster. All the other nodes sense that the expansion has occurred.



Note

If a GFS2 file system is full, the **gfs2_jadd** will fail, even if the logical volume containing the file system has been extended and is larger than the file system. This is because in a GFS2 file system, journals are plain files rather than embedded metadata, so simply extending the underlying logical volume will not provide space for the journals.

Before adding journals to a GFS file system, you can use the **journals** option of the **gfs2_tool** to find out how many journals the GFS2 file system currently contains. The following example displays the number and size of the journals in the file system mounted at **/mnt/gfs2**.

```
[root@roth-01 ../cluster/gfs2]# gfs2_tool journals /mnt/gfs2
journal2 - 128MB
journal1 - 128MB
journal0 - 128MB
3 journal(s) found.
```

Usage

```
gfs2_jadd -j Number MountPoint
```

Number

Specifies the number of new journals to be added.

MountPoint

Specifies the directory where the GFS2 file system is mounted.

Examples

In this example, one journal is added to the file system on the **/mygfs2** directory.


```
gfs2_jadd -j1 /mygfs2
```

In this example, two journals are added to the file system on the **/mygfs2** directory.

```
gfs2_jadd -j2 /mygfs2
```

Complete Usage

```
gfs2_jadd [Options] {MountPoint | Device} [MountPoint | Device]
```

MountPoint

Specifies the directory where the GFS2 file system is mounted.

Device

Specifies the device node of the file system.

[Table 3.4, “GFS2-specific Options Available When Adding Journals”](#) describes the GFS2-specific options that can be used when adding journals to a GFS2 file system.

Flag	Parameter	Description
-h		Help. Displays short usage message.
-J	<i>MegaBytes</i>	Specifies the size of the new journals in megabytes. Default journal size is 128 megabytes. The minimum size is 32 megabytes. To add journals of different sizes to the file system, the gfs2_jadd command must be run for each size journal. The size specified is rounded down so that it is a multiple of the journal-segment size that was specified when the file system was created.
-j	<i>Number</i>	Specifies the number of new journals to be added by the gfs2_jadd command. The default value is 1.
-q		Quiet. Turns down the verbosity level.
-v		Displays command version information.

Table 3.4. GFS2-specific Options Available When Adding Journals

3.8. Data Journaling

Ordinarily, GFS2 writes only metadata to its journal. File contents are subsequently written to disk by the kernel's periodic sync that flushes file system buffers. An **fsync()** call on a file causes the file's data to be written to disk immediately. The call returns when the disk reports that all data is safely written.

Data journaling can result in a reduced **fsync()** time, especially for small files, because the file data is written to the journal in addition to the metadata. An **fsync()** returns as soon as the data is written to the journal, which can be substantially faster than the time it takes to write the file data to the main file system.

Applications that rely on **fsync()** to sync file data may see improved performance by using data journaling. Data journaling can be enabled automatically for any GFS2 files created in a flagged

directory (and all its subdirectories). Existing files with zero length can also have data journaling turned on or off.

Enabling data journaling on a directory sets the directory to "inherit jdata", which indicates that all files and directories subsequently created in that directory are journaled. You can enable and disable data journaling on a file or a directory with either of the following methods:

- Executing the **chattr +j** or **chattr -j** command on the file or directory
- Setting or clearing the **jdata** flag on the file or directory with the **gfs2_tool setflag** or **gfs2_tool clearflag** command

Using the **chattr** command is the preferred way to enable and disable data journaling on a file or directory.

3.8.1. Enabling and Disabling Data Journaling with the **chattr** Command

You can enable and disable data journaling on a file with the **chattr** command. The following commands enable data journaling on the **/mnt/gfs2/gfs2_dir/newfile** file and then check whether the flag has been set properly.

```
[root@roth-01 ~]# chattr +j /mnt/gfs2/gfs2_dir/newfile
[root@roth-01 ~]# lsattr /mnt/gfs2/gfs2_dir
-----j--- /mnt/gfs2/gfs2_dir/newfile
```

The following commands disable data journaling on the **/mnt/gfs2/gfs2_dir/newfile** file and then check whether the flag has been set properly.

```
[root@roth-01 ~]# chattr -j /mnt/gfs2/gfs2_dir/newfile
[root@roth-01 ~]# lsattr /mnt/gfs2/gfs2_dir
----- /mnt/gfs2/gfs2_dir/newfile
```

You can also use the **chattr** command to set the **j** flag on a directory. When you set this flag for a directory, all files and directories subsequently created in that directory are journaled. The following set of commands sets the **j** flag on the **gfs2_dir** directory, then checks whether the flag has been set properly. After this, the commands create a new file called **newfile** in the **/mnt/gfs2/gfs2_dir** directory and then check whether the **j** flag has been set for the file. Since the **j** flag is set for the directory, then **newfile** should also have journaling enabled.

```
[root@roth-01 ~]# chattr +j /mnt/gfs2/gfs2_dir
[root@roth-01 ~]# lsattr /mnt/gfs2
-----j--- /mnt/gfs2/gfs2_dir
[root@roth-01 ~]# touch /mnt/gfs2/gfs2_dir/newfile
[root@roth-01 ~]# lsattr /mnt/gfs2/gfs2_dir
-----j--- /mnt/gfs2/gfs2_dir/newfile
```

3.8.2. Enabling and Disabling Data Journaling with the **gfs2_tool** Command

The following command enables data journaling on the GFS2 file **/mnt/gfs2/gfs2file**

```
[root@roth-01 ~]# gfs2_tool setflag jdata /mnt/gfs2/gfs2file
```

You can use the **lsattr** command to verify that the **jdata** flag has been set.

```
[root@roth-01 ~]# lsattr /mnt/gfs2
-----j--- /mnt/gfs2/gfs2file
```

When the **jdata** flag is set for a directory, all files and directories subsequently created in that directory are journaled.

The following commands create a new subdirectory called **gfs2_dir** in the **/mnt/gfs2** directory, set the **jdata** flag for the directory, then check whether the flag has been set properly. After this, The commands create the file **gfs2file** in the directory then verify that the flag for that file has been set properly.

```
[root@roth-01 ~]# mkdir /mnt/gfs2/gfs2_dir
[root@roth-01 ~]# gfs2_tool setflag jdata /mnt/gfs2/gfs2_dir
[root@roth-01 ~]# lsattr /mnt/gfs2
-----j--- /mnt/gfs2/gfs2_dir
[root@roth-01 ~]# touch /mnt/gfs2/gfs2_dir/newfile
[root@roth-01 ~]# lsattr /mnt/gfs2/gfs2_dir
-----j--- /mnt/gfs2/gfs2_dir/newfile
```

To clear the **jdata** flag from a file or directory, you can use the **gfs2_tool clearflag jdata** command, as in the following example.

```
[root@roth-01 ~]# gfs2_tool clearflag jdata /mnt/gfs2/gfs2file
[root@roth-01 ~]# lsattr /mnt/gfs2
----- /mnt/gfs2/gfs2file
```

3.9. Configuring **atime** Updates

Each file inode and directory inode has three time stamps associated with it:

- **ctime** — The last time the inode status was changed
- **mtime** — The last time the file (or directory) data was modified
- **atime** — The last time the file (or directory) data was accessed

If **atime** updates are enabled as they are by default on GFS2 and other Linux file systems then every time a file is read, its inode needs to be updated.

Because few applications use the information provided by **atime**, those updates can require a significant amount of unnecessary write traffic and file locking traffic. That traffic can degrade performance; therefore, it may be preferable to turn off or reduce the frequency of **atime** updates.

Two methods of reducing the effects of **atime** updating are available:

- Mount with **relatime** (relative atime), which updates the **atime** if the previous **atime** update is older than the **mtime** or **ctime** update.

- Mount with **noatime**, which disables **atime** updates on that file system.

3.9.1. Mount with **relatime**

The **relatime** (relative atime) Linux mount option can be specified when the file system is mounted. This specifies that the **atime** is updated if the previous **atime** update is older than the **mtime** or **ctime** update.

Usage

```
mount BlockDevice MountPoint -o relatime
```

BlockDevice

Specifies the block device where the GFS2 file system resides.

MountPoint

Specifies the directory where the GFS2 file system should be mounted.

Example

In this example, the GFS2 file system resides on the **/dev/vg01/lvol10** and is mounted on directory **/mygfs2**. The **atime** updates take place only if the previous **atime** update is older than the **mtime** or **ctime** update.

```
mount /dev/vg01/lvol10 /mygfs2 -o relatime
```

3.9.2. Mount with **noatime**

The **noatime** Linux mount option can be specified when the file system is mounted, which disables **atime** updates on that file system.

Usage

```
mount BlockDevice MountPoint -o noatime
```

BlockDevice

Specifies the block device where the GFS2 file system resides.

MountPoint

Specifies the directory where the GFS2 file system should be mounted.

Example

In this example, the GFS2 file system resides on the **/dev/vg01/lvol10** and is mounted on directory **/mygfs2** with **atime** updates turned off.

```
mount /dev/vg01/lvol10 /mygfs2 -o noatime
```

3.10. Suspending Activity on a File System

You can suspend write activity to a file system by using the **gfs2_tool freeze** command. Suspending write activity allows hardware-based device snapshots to be used to capture the file system in a consistent state. The **gfs2_tool unfreeze** command ends the suspension.

Usage

Start Suspension

```
gfs2_tool freeze MountPoint
```

End Suspension

```
gfs2_tool unfreeze MountPoint
```

MountPoint

Specifies the file system.

Examples

This example suspends writes to file system **/mygfs2**.

```
gfs2_tool freeze /mygfs2
```

This example ends suspension of writes to file system **/mygfs2**.

```
gfs2_tool unfreeze /mygfs2
```

3.11. Repairing a File System

When nodes fail with the file system mounted, file system journaling allows fast recovery. However, if a storage device loses power or is physically disconnected, file system corruption may occur. (Journaling cannot be used to recover from storage subsystem failures.) When that type of corruption occurs, you can recover the GFS2 file system by using the **fsck.gfs2** command.



Warning

The **fsck.gfs2** command must be run only on a file system that is unmounted from all nodes.



Note

If you have previous experience using the **gfs_fsck** command on GFS file systems, note that the **fsck.gfs2** command differs from some earlier releases of **gfs_fsck** in the in the following ways:

- You cannot set the interactive mode with **Ctrl+C**. Pressing **Ctrl+C** cancels the **fsck.gfs2** command. Do *not* press **Ctrl+C** unless you want to cancel the command.
- You can increase the level of verbosity by using the **-v** flag. Adding a second **-v** flag increases the level again.
- You can decrease the level of verbosity by using the **-q** flag. Adding a second **-q** flag decreases the level again.
- The **-n** option opens a file system as read-only and answers **no** to any queries automatically. The option provides a way of trying the command to reveal errors without actually allowing the **fsck.gfs2** command to take effect.

Refer to the **gfs2.fsck** man page for additional information about other command options.

Usage

```
fsck.gfs2 -y BlockDevice
```

-y

The **-y** flag causes all questions to be answered with **yes**. With the **-y** flag specified, the **fsck.gfs2** command does not prompt you for an answer before making changes.

BlockDevice

Specifies the block device where the GFS2 file system resides.

Example

In this example, the GFS2 file system residing on block device **/dev/testvol/testlv** is repaired. All queries to repair are automatically answered with **yes**.

```
[root@dash-01 ~]# fsck.gfs2 -y /dev/testvg/testlv
Initializing fsck
Validating Resource Group index.
Level 1 RG check.
(level 1 passed)
Clearing journals (this may take a while)...
Journals cleared.
Starting pass1
Pass1 complete
Starting pass1b
Pass1b complete
Starting pass1c
Pass1c complete
Starting pass2
Pass2 complete
Starting pass3
Pass3 complete
Starting pass4
Pass4 complete
Starting pass5
```

```
Pass5 complete
Writing changes to disk
fsck.gfs2 complete
```

3.12. Bind Mounts and Context-Dependent Path Names

GFS2 file systems do not provide support for Context-Dependent Path Names (CDPNs), which allow you to create symbolic links that point to variable destination files or directories. For this functionality in GFS2, you can use the **bind** option of the **mount** command.

The **bind** option of the **mount** command allows you to remount part of a file hierarchy at a different location while it is still available at the original location. The format of this command is as follows.

```
mount --bind olddir newdir
```

After executing this command, the contents of the *olddir* directory are available at two locations: *olddir* and *newdir*. You can also use this option to make an individual file available at two locations.

For example, after executing the following commands the contents of **/root/tmp** will be identical to the contents of the previously mounted **/var/log** directory.

```
[root@menscryfa ~]# cd ~root
[root@menscryfa ~]# mkdir ./tmp
[root@menscryfa ~]# mount --bind /var/log /tmp
```

Alternately, you can use an entry in the **/etc/fstab** file to achieve the same results at mount time. The following **/etc/fstab** entry will result in the contents of **/root/tmp** being identical to the contents of the **/var/log** directory.

/var/log	/root/tmp	none	bind	0 0
----------	-----------	------	------	-----

After you have mounted the file system, you can use the **mount** command to see that the file system has been mounted, as in the following example.

```
[root@menscryfa ~]# mount | grep /tmp
/var/log on /root/tmp type none (rw,bind)
```

With a file system that supports Context-Dependent Path Names, you might have defined the **/bin** directory as a Context-Dependent Path Name that would resolve to one of the following paths, depending on the system architecture.

```
/usr/i386-bin
/usr/x86_64-bin
/usr/ppc64-bin
```

You can achieve this same functionality by creating an empty **/bin** directory. Then, using a script or an entry in the **/etc/fstab** file, you can mount each of the individual architecture directories onto the

/bin directory with a **mount -bind** command. For example, you can use the following command as a line in a script.

```
mount --bind /usr/i386-bin /bin
```

Alternately, you can use the following entry in the **/etc/fstab** file.

/usr/i386-bin	/bin	none	bind	0 0
---------------	------	------	------	-----

A bind mount can provide greater flexibility than a Context-Dependent Path Name, since you can use this feature to mount different directories according to any criteria you define (such as the value of **%fill** for the file system). Context-Dependent Path Names are more limited in what they can encompass. Note, however, that you will need to write your own script to mount according to a criteria such as the value of **%fill**.



Warning

When you mount a file system with the **bind** option and the original file system was mounted **rw**, the new file system will also be mounted **rw** even if you use the **ro** flag; the **ro** flag is silently ignored. In this case, the new file system might be marked as **ro** in the **/proc/mounts** directory, which may be misleading.

3.13. Bind Mounts and File System Mount Order

When you use the **bind** option of the **mount** command, you must be sure that the file systems are mounted in the correct order. In the following example, the **/var/log** directory must be mounted before executing the bind mount on the **/tmp** directory:

```
# mount --bind /var/log /tmp
```

The ordering of file system mounts is determined as follows:

- In general, file system mount order is determined by the order in which the file systems appear in the **fstab** file. The exceptions to this ordering are file systems mounted with the **_netdev** flag or filesystems that have their own **init** scripts.
- A file system with its own **init** script is mounted later in the initialization process, after the file systems in the **fstab** file.
- File systems mounted with the **_netdev** flag are mounted when the network has been enabled on the system.

If your configuration requires that you create a bind mount on which to mount a GFS2 file system, you can order your **fstab** file as follows:

1. Mount local filesystems that are required for the bind mount.
2. Bind mount the directory on which to mount the GFS2 file system.
3. Mount the GFS2 file system.

If your configuration requires that you bind mount a local directory or file system onto a GFS2 file system, listing the file systems in the correct order in the **fstab** file will not mount the file systems correctly since the GFS2 file system will not be mounted until the GFS2 **init** script is run. In this case, you should write an **init** script to execute the bind mount so that the bind mount will not take place until after the GFS2 file system is mounted.

The following script is an example of a custom **init** script. This script performs a bind mount of two directories onto two directories of a GFS2 filesystem. In this example, there is an existing GFS2 mount point at **/mnt/gfs2a**, which is mounted when the GFS2 **init** script runs, after cluster startup.

In this example script, the values of the **chkconfig** statement indicate the following:

- 345 indicates the run levels that the script will be started in
- 29 is the start priority, which in this case indicates that the script will run at startup time after the GFS2 **init** script, which has a start priority of 26
- 73 is the stop priority, which in this case indicates that the script will be stopped during shutdown before the GFS2 script, which has a stop priority of 74

The start and stop values indicate that you can manually perform the indicated action by executing a **service start** and a **service stop** command. For example, if the script is named **fredwilma**, then you can execute **service fredwilma start**.

This script should be put in the **/etc/init.d** directory with the same permissions as the other scripts in that directory. You can then execute a **chkconfig on** command to link the script to the indicated run levels. For example, if the script is named **fredwilma**, then you can execute **chkconfig fredwilma on**.

```
#!/bin/bash
#
# chkconfig: 345 29 73
# description: mount/unmount my custom bind mounts onto a gfs2 subdirectory
#
### BEGIN INIT INFO
# Provides:
### END INIT INFO

. /etc/init.d/functions
case "$1" in
    start)
        # In this example, fred and wilma want their home directories
        # bind-mounted over the gfs2 directory /mnt/gfs2a, which has
        # been mounted as /mnt/gfs2a
        mkdir -p /mnt/gfs2a/home/fred &> /dev/null
        mkdir -p /mnt/gfs2a/home/wilma &> /dev/null
        /bin/mount --bind /mnt/gfs2a/home/fred /home/fred
        /bin/mount --bind /mnt/gfs2a/home/wilma /home/wilma
        ;;

    stop)
        /bin/umount /mnt/gfs2a/home/fred
        /bin/umount /mnt/gfs2a/home/wilma
        ;;

    status)

```

```
;;
restart)
    $0 stop
    $0 start
    ;;

reload)
    $0 start
    ;;
*)
    echo $"Usage: $0 {start|stop|restart|reload|status}"
    exit 1
esac
exit 0
```

3.14. The GFS2 Withdraw Function

The GFS2 *withdraw* function is a data integrity feature of GFS2 file systems in a cluster. If the GFS2 kernel module detects an inconsistency in a GFS2 file system following an I/O operation, the node will withdraw itself from the cluster. This stops the I/O operation and waits for further I/O operations to error out, preventing further damage. You can then unmount the file system, perform file system recovery with an **fsck** command, and reboot the node. You can also stop any services or applications manually before you reboot. The GFS2 withdraw function is less severe than a kernel panic, which would cause another node to fence the node.

An example of an inconsistency that would yield a GFS2 withdraw is an incorrect block count. When the GFS2 kernel deletes a file from a file system, it systematically removes all the data and metadata blocks associated with that file. When it is done, it checks the block count. If the block count is not one (meaning all that is left is the disk inode itself), that indicates a file system inconsistency since the block count did not match the list of blocks found.

Appendix A. Converting a File System from GFS to GFS2

Use the **gfs2_convert** command. To convert a GFS file system to a GFS2 file system.

1. Unmount the GFS file system from all nodes in the cluster.
2. Execute the **gfs_fsck** command on the GFS file system to ensure there is no file system corruption.
3. Remount and make a backup of your existing GFS file system.
4. Unmount the file system and execute **gfs2_convert *gfsfilesystem***. The system will display warnings and confirmation questions before converting *gfsfilesystem* to GFS2.

The following example converts a GFS filesystem on block device **/dev/testvg/testlv** to a GFS2 filesystem.

```
[root@dash-01 ~]# gfs2_convert /dev/testvg/testlv
gfs2_convert version 2 (built Sep 25 2007 12:41:29)
Copyright (C) Red Hat, Inc. 2004-2006 All rights reserved.

This program will convert a gfs1 filesystem to a gfs2 filesystem.
WARNING: This can't be undone. It is strongly advised that you:

    1. Back up your entire filesystem first.
    2. Run gfs_fsck first to ensure filesystem integrity.
    3. Make sure the filesystem is NOT mounted from any node.
    4. Make sure you have the latest software versions.
Convert /dev/testvg/testlv from GFS1 to GFS2? (y/n)y
Converting resource groups.
Converting inodes.
6 inodes converted.
Fixing file and directory information.
1 directories, 2 dirents fixed.
Converting journals.
Converting journal space to rg space.
Building system structures.
Removing obsolete gfs1 structures.
Committing changes to disk.
/dev/gfs2testvg/gfs2testlv: filesystem converted successfully to gfs2.
[root@dash-01 ~]#
```

Appendix B. Revision History

Revision 5.4-2 Wed Feb 03 2010

Steven Levine slevine@redhat.com

Resolves: #561560

Adds note indicating that GFS2 is not supported in clusters greater than 16 nodes.

Revision 5.4-1 Tue Aug 18 2009

Steven Levine slevine@redhat.com

Resolves: #515807

Adds note clarifying that you cannot reduce the size of an existing file system.

Resolves: #480002

Adds caveat about unmounting a file system manually if you mounted it manually, adds section on bind mounting a non-GFS2 file system to a GFS2 file system, adds sample custom init script.

Resolves: #458604

Adds section on GFS2 withdraw function.

Resolves: #498292

Clarifies documentation on adding journals.

Revision 1.0 Thu Jan 29 2009

Index

A

- acl mount option, 12
- adding journals to a file system, 22
- atime, configuring updates, 25
 - mounting with noatime, 26
 - mounting with relatime, 26
- audience, v

B

- bind mount
 - mount order, 30
- bind mounts, 29

C

- configuration, before, 2
- configuration, initial, 7
 - prerequisite tasks, 7

D

- data journaling, 23

F

- feedback, vi, vi
- file system
 - adding journals, 22
 - atime, configuring updates, 25
 - mounting with noatime, 26
 - mounting with relatime, 26
 - bind mounts, 29
 - context-dependent path names (CDPNs), 29
 - data journaling, 23
 - growing, 20
 - making, 9
 - mount order, 30
 - mounting, 12, 15
 - quota management, 15
 - displaying quota limits, 17
 - enabling quota accounting, 20
 - enabling/disabling quota enforcement, 19
 - setting quotas, 15
 - synchronizing quotas, 18
 - repairing, 27
 - suspending activity, 27
 - unmounting, 14, 15
- fsck.gfs2 command, 27

G

- GFS2
 - atime, configuring updates, 25
 - mounting with noatime, 26
 - mounting with relatime, 26
 - managing, 9
 - quota management, 15
 - displaying quota limits, 17
 - enabling quota accounting, 20
 - enabling/disabling quota enforcement, 19
 - setting quotas, 15
 - synchronizing quotas, 18
 - withdraw function, 32
- GFS2 file system maximum size, 1
- GFS2-specific options for adding journals table, 23
- GFS2-specific options for expanding file systems table, 21
- gfs2_grow command, 20
- gfs2_jadd command, 22
- gfs2_quota command, 15
- growing a file system, 20

I

- initial tasks
 - setup, initial, 7
- introduction, v
 - audience, v

M

- making a file system, 9
- managing GFS2, 9
- maximum size, GFS2 file system, 1
- mkfs command, 9
- mkfs.gfs2 command options table, 11
- mount command, 12
- mount table, 13
- mounting a file system, 12, 15

O

- overview, 1
 - configuration, before, 2

P

- path names, context-dependent (CDPNs), 29
- preface (see introduction)
- prerequisite tasks
 - configuration, initial, 7

Q

- quota management, 15
 - displaying quota limits, 17
 - enabling quota accounting, 20
 - enabling/disabling quota enforcement, 19
 - setting quotas, 15
 - synchronizing quotas, 18
- quota= mount option, 16
- quota_quantum tunable parameter, 18

R

- repairing a file system, 27

S

- setup, initial
 - initial tasks, 7
- suspending activity on a file system, 27
- system hang at unmount, 15

T

- tables
 - GFS2-specific options for adding journals, 23
 - GFS2-specific options for expanding file systems, 21
 - mkfs.gfs2 command options, 11
 - mount options, 13

U

- umount command, 14
- unmount, system hang, 15
- unmounting a file system, 14, 15

W

- withdraw function, GFS2, 32