



Red Hat Enterprise Linux 9.0

Managing file systems

Creating, modifying, and administering file systems in Red Hat Enterprise Linux 9

Red Hat Enterprise Linux 9.0 Managing file systems

Creating, modifying, and administering file systems in Red Hat Enterprise Linux 9

Legal Notice

Copyright © 2023 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, the Red Hat logo, JBoss, OpenShift, 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.

Java[®] is a registered trademark of Oracle and/or its affiliates.

XFS[®] is a trademark of Silicon Graphics International Corp. or its subsidiaries in the United States and/or other countries.

MySQL[®] is a registered trademark of MySQL AB in the United States, the European Union and other countries.

Node.js[®] is an official trademark of Joyent. Red Hat is not formally related to or endorsed by the official Joyent Node.js open source or commercial project.

The OpenStack[®] Word Mark and OpenStack logo are either registered trademarks/service marks or trademarks/service marks of the OpenStack Foundation, in the United States and other countries and are used with the OpenStack Foundation's permission. We are not affiliated with, endorsed or sponsored by the OpenStack Foundation, or the OpenStack community.

All other trademarks are the property of their respective owners.

Abstract

Red Hat Enterprise Linux supports a variety of file systems. Each type of file system solves different problems and their usage is application specific. Use the information about the key differences and considerations to select and deploy the appropriate file system based on your specific application requirements. The supported file systems include local on-disk file systems XFS and ext4, network and client-and-server file systems NFS and SMB, as well as a combined local storage and file system management solution, Stratis. You can perform several operations with a file system such as creating, mounting, backing up, restoring, checking and repairing, as well as limiting the storage space by using quotas.

Table of Contents

MAKING OPEN SOURCE MORE INCLUSIVE	8
PROVIDING FEEDBACK ON RED HAT DOCUMENTATION	9
CHAPTER 1. OVERVIEW OF AVAILABLE FILE SYSTEMS	10
1.1. TYPES OF FILE SYSTEMS	10
1.2. LOCAL FILE SYSTEMS	11
1.3. THE XFS FILE SYSTEM	11
1.4. THE EXT4 FILE SYSTEM	12
1.5. COMPARISON OF XFS AND EXT4	13
1.6. CHOOSING A LOCAL FILE SYSTEM	14
1.7. NETWORK FILE SYSTEMS	15
1.8. SHARED STORAGE FILE SYSTEMS	15
1.9. CHOOSING BETWEEN NETWORK AND SHARED STORAGE FILE SYSTEMS	16
1.10. VOLUME-MANAGING FILE SYSTEMS	17
CHAPTER 2. MANAGING LOCAL STORAGE USING RHEL SYSTEM ROLES	18
2.1. INTRODUCTION TO THE STORAGE RHEL SYSTEM ROLE	18
2.2. PARAMETERS THAT IDENTIFY A STORAGE DEVICE IN THE STORAGE RHEL SYSTEM ROLE	18
2.3. EXAMPLE ANSIBLE PLAYBOOK TO CREATE AN XFS FILE SYSTEM ON A BLOCK DEVICE	19
2.4. EXAMPLE ANSIBLE PLAYBOOK TO PERSISTENTLY MOUNT A FILE SYSTEM	20
2.5. EXAMPLE ANSIBLE PLAYBOOK TO MANAGE LOGICAL VOLUMES	20
2.6. EXAMPLE ANSIBLE PLAYBOOK TO ENABLE ONLINE BLOCK DISCARD	21
2.7. EXAMPLE ANSIBLE PLAYBOOK TO CREATE AND MOUNT AN EXT4 FILE SYSTEM	22
2.8. EXAMPLE ANSIBLE PLAYBOOK TO CREATE AND MOUNT AN EXT3 FILE SYSTEM	22
2.9. EXAMPLE ANSIBLE PLAYBOOK TO RESIZE AN EXISTING EXT4 OR EXT3 FILE SYSTEM USING THE STORAGE RHEL SYSTEM ROLE	23
2.10. EXAMPLE ANSIBLE PLAYBOOK TO RESIZE AN EXISTING FILE SYSTEM ON LVM USING THE STORAGE RHEL SYSTEM ROLE	24
2.11. EXAMPLE ANSIBLE PLAYBOOK TO CREATE A SWAP VOLUME USING THE STORAGE RHEL SYSTEM ROLE	25
2.12. CONFIGURING A RAID VOLUME USING THE STORAGE SYSTEM ROLE	26
2.13. CONFIGURING AN LVM POOL WITH RAID USING THE STORAGE RHEL SYSTEM ROLE	27
2.14. EXAMPLE ANSIBLE PLAYBOOK TO COMPRESS AND DEDUPLICATE A VDO VOLUME ON LVM USING THE STORAGE RHEL SYSTEM ROLE	28
2.15. CREATING A LUKS2 ENCRYPTED VOLUME USING THE STORAGE RHEL SYSTEM ROLE	29
2.16. EXAMPLE ANSIBLE PLAYBOOK TO EXPRESS POOL VOLUME SIZES AS PERCENTAGE USING THE STORAGE RHEL SYSTEM ROLE	31
2.17. ADDITIONAL RESOURCES	32
CHAPTER 3. MOUNTING NFS SHARES	33
3.1. INTRODUCTION TO NFS	33
3.2. SUPPORTED NFS VERSIONS	33
Default NFS version	33
Features of minor NFS versions	33
3.3. SERVICES REQUIRED BY NFS	34
The RPC services with NFSv4	35
3.4. NFS HOST NAME FORMATS	35
3.5. INSTALLING NFS	36
3.6. DISCOVERING NFS EXPORTS	36
3.7. MOUNTING AN NFS SHARE WITH MOUNT	36
3.8. COMMON NFS MOUNT OPTIONS	37
3.9. STORING USER SETTINGS OVER NFS	39

3.10. ADDITIONAL RESOURCES	39
CHAPTER 4. EXPORTING NFS SHARES	40
4.1. INTRODUCTION TO NFS	40
4.2. SUPPORTED NFS VERSIONS	40
Default NFS version	40
Features of minor NFS versions	40
4.3. THE TCP AND UDP PROTOCOLS IN NFSV3 AND NFSV4	41
4.4. SERVICES REQUIRED BY NFS	41
The RPC services with NFSv4	42
4.5. NFS HOST NAME FORMATS	42
4.6. NFS SERVER CONFIGURATION	43
4.6.1. The /etc/exports configuration file	43
Export entry	43
Default options	44
Default and overridden options	45
4.6.2. The exportfs utility	45
Common exportfs options	45
4.7. NFS AND RPCBIND	46
4.8. INSTALLING NFS	46
4.9. STARTING THE NFS SERVER	46
4.10. TROUBLESHOOTING NFS AND RPCBIND	47
4.11. CONFIGURING THE NFS SERVER TO RUN BEHIND A FIREWALL	48
4.11.1. Configuring the NFSv3-enabled server to run behind a firewall	48
4.11.2. Configuring the NFSv4-only server to run behind a firewall	49
4.11.3. Configuring an NFSv3 client to run behind a firewall	50
4.11.4. Configuring an NFSv4 client to run behind a firewall	51
4.12. EXPORTING RPC QUOTA THROUGH A FIREWALL	51
4.13. ENABLING NFS OVER RDMA (NFSORDMA)	52
4.14. ADDITIONAL RESOURCES	53
CHAPTER 5. CONFIGURING AN NFSV4-ONLY SERVER	54
5.1. BENEFITS AND DRAWBACKS OF AN NFSV4-ONLY SERVER	54
5.2. CONFIGURING THE NFS SERVER TO SUPPORT ONLY NFSV4	54
5.3. VERIFYING THE NFSV4-ONLY CONFIGURATION	55
CHAPTER 6. SECURING NFS	57
6.1. NFS SECURITY WITH AUTH_SYS AND EXPORT CONTROLS	57
6.2. NFS SECURITY WITH AUTH_GSS	57
6.3. CONFIGURING AN NFS SERVER AND CLIENT TO USE KERBEROS	57
6.4. NFSV4 SECURITY OPTIONS	58
6.5. FILE PERMISSIONS ON MOUNTED NFS EXPORTS	59
CHAPTER 7. ENABLING PNFS SCSI LAYOUTS IN NFS	60
7.1. THE PNFS TECHNOLOGY	60
7.2. PNFS SCSI LAYOUTS	60
Operations between the client and the server	60
Device reservations	60
7.3. CHECKING FOR A SCSI DEVICE COMPATIBLE WITH PNFS	61
7.4. SETTING UP PNFS SCSI ON THE SERVER	62
7.5. SETTING UP PNFS SCSI ON THE CLIENT	62
7.6. RELEASING THE PNFS SCSI RESERVATION ON THE SERVER	63
CHAPTER 8. MONITORING PNFS SCSI LAYOUTS FUNCTIONALITY	65

8.1. CHECKING PNFS SCSI OPERATIONS FROM THE SERVER USING NFSSTAT	65
8.2. CHECKING PNFS SCSI OPERATIONS FROM THE CLIENT USING MOUNTSTATS	65
CHAPTER 9. GETTING STARTED WITH FS-CACHE	67
9.1. OVERVIEW OF THE FS-CACHE	67
9.2. PERFORMANCE GUARANTEE	68
9.3. SETTING UP A CACHE	68
9.4. CACHE CULL LIMITS CONFIGURATION	69
9.5. RETRIEVING STATISTICAL INFORMATION FROM THE FSCACHE KERNEL MODULE	70
9.6. FS-CACHE REFERENCES	71
CHAPTER 10. USING THE CACHE WITH NFS	72
10.1. CONFIGURING NFS CACHE SHARING	72
10.2. CACHE LIMITATIONS WITH NFS	73
CHAPTER 11. MOUNTING AN SMB SHARE ON RED HAT ENTERPRISE LINUX	75
11.1. SUPPORTED SMB PROTOCOL VERSIONS	75
11.2. UNIX EXTENSIONS SUPPORT	76
11.3. MANUALLY MOUNTING AN SMB SHARE	76
11.4. MOUNTING AN SMB SHARE AUTOMATICALLY WHEN THE SYSTEM BOOTS	77
11.5. AUTHENTICATING TO AN SMB SHARE USING A CREDENTIALS FILE	77
11.6. FREQUENTLY USED MOUNT OPTIONS	78
CHAPTER 12. PERFORMING A MULTI-USER SMB MOUNT	80
12.1. MOUNTING A SHARE WITH THE MULTIUSER OPTION	80
12.2. VERIFYING IF AN SMB SHARE IS MOUNTED WITH THE MULTIUSER OPTION	81
12.3. ACCESSING A SHARE AS A USER	81
CHAPTER 13. OVERVIEW OF PERSISTENT NAMING ATTRIBUTES	82
13.1. DISADVANTAGES OF NON-PERSISTENT NAMING ATTRIBUTES	82
13.2. FILE SYSTEM AND DEVICE IDENTIFIERS	82
File system identifiers	83
Device identifiers	83
Recommendations	83
13.3. DEVICE NAMES MANAGED BY THE UDEV MECHANISM IN /DEV/DISK/	83
13.3.1. File system identifiers	83
The UUID attribute in /dev/disk/by-uuid/	83
The Label attribute in /dev/disk/by-label/	84
13.3.2. Device identifiers	84
The WWID attribute in /dev/disk/by-id/	84
The Partition UUID attribute in /dev/disk/by-partuuid	85
The Path attribute in /dev/disk/by-path/	85
13.4. THE WORLD WIDE IDENTIFIER WITH DM MULTIPATH	85
13.5. LIMITATIONS OF THE UDEV DEVICE NAMING CONVENTION	86
13.6. LISTING PERSISTENT NAMING ATTRIBUTES	87
13.7. MODIFYING PERSISTENT NAMING ATTRIBUTES	88
CHAPTER 14. PARTITION OPERATIONS WITH PARTED	89
14.1. VIEWING THE PARTITION TABLE WITH PARTED	89
14.2. CREATING A PARTITION TABLE ON A DISK WITH PARTED	90
14.3. CREATING A PARTITION WITH PARTED	91
14.4. REMOVING A PARTITION WITH PARTED	92
14.5. RESIZING A PARTITION WITH PARTED	94
CHAPTER 15. STRATEGIES FOR REPARTITIONING A DISK	96

15.1. USING UNPARTITIONED FREE SPACE	96
15.2. USING SPACE FROM AN UNUSED PARTITION	96
15.3. USING FREE SPACE FROM AN ACTIVE PARTITION	97
15.3.1. Destructive repartitioning	97
15.3.2. Non-destructive repartitioning	98
CHAPTER 16. GETTING STARTED WITH XFS	101
16.1. THE XFS FILE SYSTEM	101
16.2. COMPARISON OF TOOLS USED WITH EXT4 AND XFS	102
CHAPTER 17. CREATING AN XFS FILE SYSTEM	103
17.1. CREATING AN XFS FILE SYSTEM WITH MKFS.XFS	103
CHAPTER 18. BACKING UP AN XFS FILE SYSTEM	104
18.1. FEATURES OF XFS BACKUP	104
18.2. BACKING UP AN XFS FILE SYSTEM WITH XFS_DUMP	104
18.3. ADDITIONAL RESOURCES	105
CHAPTER 19. RESTORING AN XFS FILE SYSTEM FROM BACKUP	106
19.1. FEATURES OF RESTORING XFS FROM BACKUP	106
19.2. RESTORING AN XFS FILE SYSTEM FROM BACKUP WITH XFSRESTORE	106
19.3. INFORMATIONAL MESSAGES WHEN RESTORING AN XFS BACKUP FROM A TAPE	107
19.4. ADDITIONAL RESOURCES	108
CHAPTER 20. INCREASING THE SIZE OF AN XFS FILE SYSTEM	109
20.1. INCREASING THE SIZE OF AN XFS FILE SYSTEM WITH XFS_GROWFS	109
CHAPTER 21. CONFIGURING XFS ERROR BEHAVIOR	110
21.1. CONFIGURABLE ERROR HANDLING IN XFS	110
21.2. CONFIGURATION FILES FOR SPECIFIC AND UNDEFINED XFS ERROR CONDITIONS	110
21.3. SETTING XFS BEHAVIOR FOR SPECIFIC CONDITIONS	111
21.4. SETTING XFS BEHAVIOR FOR UNDEFINED CONDITIONS	111
21.5. SETTING THE XFS UNMOUNT BEHAVIOR	112
CHAPTER 22. CHECKING AND REPAIRING A FILE SYSTEM	113
22.1. SCENARIOS THAT REQUIRE A FILE SYSTEM CHECK	113
22.2. POTENTIAL SIDE EFFECTS OF RUNNING FSCK	114
22.3. ERROR-HANDLING MECHANISMS IN XFS	114
Unclean unmounts	114
Corruption	114
22.4. CHECKING AN XFS FILE SYSTEM WITH XFS_REPAIR	115
22.5. REPAIRING AN XFS FILE SYSTEM WITH XFS_REPAIR	116
22.6. ERROR HANDLING MECHANISMS IN EXT2, EXT3, AND EXT4	117
22.7. CHECKING AN EXT2, EXT3, OR EXT4 FILE SYSTEM WITH E2FSCK	117
22.8. REPAIRING AN EXT2, EXT3, OR EXT4 FILE SYSTEM WITH E2FSCK	118
CHAPTER 23. MOUNTING FILE SYSTEMS	119
23.1. THE LINUX MOUNT MECHANISM	119
23.2. LISTING CURRENTLY MOUNTED FILE SYSTEMS	119
23.3. MOUNTING A FILE SYSTEM WITH MOUNT	120
23.4. MOVING A MOUNT POINT	121
23.5. UNMOUNTING A FILE SYSTEM WITH UMount	121
23.6. COMMON MOUNT OPTIONS	122
CHAPTER 24. SHARING A MOUNT ON MULTIPLE MOUNT POINTS	124
24.1. TYPES OF SHARED MOUNTS	124

24.2. CREATING A PRIVATE MOUNT POINT DUPLICATE	124
24.3. CREATING A SHARED MOUNT POINT DUPLICATE	125
24.4. CREATING A SLAVE MOUNT POINT DUPLICATE	127
24.5. PREVENTING A MOUNT POINT FROM BEING DUPLICATED	128
CHAPTER 25. PERSISTENTLY MOUNTING FILE SYSTEMS	129
25.1. THE /ETC/FSTAB FILE	129
25.2. ADDING A FILE SYSTEM TO /ETC/FSTAB	130
CHAPTER 26. MOUNTING FILE SYSTEMS ON DEMAND	131
26.1. THE AUTOFS SERVICE	131
26.2. THE AUTOFS CONFIGURATION FILES	131
26.3. CONFIGURING AUTOFS MOUNT POINTS	133
26.4. AUTOMOUNTING NFS SERVER USER HOME DIRECTORIES WITH AUTOFS SERVICE	134
26.5. OVERRIDING OR AUGMENTING AUTOFS SITE CONFIGURATION FILES	134
26.6. USING LDAP TO STORE AUTOMOUNTER MAPS	136
26.7. USING SYSTEMD.AUTOMOUNT TO MOUNT A FILE SYSTEM ON DEMAND WITH /ETC/FSTAB	137
26.8. USING SYSTEMD.AUTOMOUNT TO MOUNT A FILE SYSTEM ON DEMAND WITH A MOUNT UNIT	138
CHAPTER 27. USING SSSD COMPONENT FROM IDM TO CACHE THE AUTOFS MAPS	140
27.1. CONFIGURING AUTOFS MANUALLY TO USE IDM SERVER AS AN LDAP SERVER	140
27.2. CONFIGURING SSSD TO CACHE AUTOFS MAPS	141
CHAPTER 28. SETTING READ-ONLY PERMISSIONS FOR THE ROOT FILE SYSTEM	143
28.1. FILES AND DIRECTORIES THAT ALWAYS RETAIN WRITE PERMISSIONS	143
28.2. CONFIGURING THE ROOT FILE SYSTEM TO MOUNT WITH READ-ONLY PERMISSIONS ON BOOT	144
CHAPTER 29. LIMITING STORAGE SPACE USAGE ON XFS WITH QUOTAS	145
29.1. DISK QUOTAS	145
29.2. THE XFS_QUOTA TOOL	145
29.3. FILE SYSTEM QUOTA MANAGEMENT IN XFS	145
29.4. ENABLING DISK QUOTAS FOR XFS	146
29.5. REPORTING XFS USAGE	146
29.6. MODIFYING XFS QUOTA LIMITS	147
29.7. SETTING PROJECT LIMITS FOR XFS	148
CHAPTER 30. LIMITING STORAGE SPACE USAGE ON EXT4 WITH QUOTAS	150
30.1. INSTALLING THE QUOTA TOOL	150
30.2. ENABLING QUOTA FEATURE ON FILE SYSTEM CREATION	150
30.3. ENABLING QUOTA FEATURE ON EXISTING FILE SYSTEMS	150
30.4. ENABLING QUOTA ENFORCEMENT	151
30.5. ASSIGNING QUOTAS PER USER	152
30.6. ASSIGNING QUOTAS PER GROUP	153
30.7. ASSIGNING QUOTAS PER PROJECT	154
30.8. SETTING THE GRACE PERIOD FOR SOFT LIMITS	155
30.9. TURNING FILE SYSTEM QUOTAS OFF	155
30.10. REPORTING ON DISK QUOTAS	156
CHAPTER 31. DISCARDING UNUSED BLOCKS	157
Requirements	157
31.1. TYPES OF BLOCK DISCARD OPERATIONS	157
Recommendations	157
31.2. PERFORMING BATCH BLOCK DISCARD	157
31.3. ENABLING ONLINE BLOCK DISCARD	158
31.4. ENABLING PERIODIC BLOCK DISCARD	158

CHAPTER 32. SETTING UP STRATIS FILE SYSTEMS	160
32.1. WHAT IS STRATIS	160
32.2. COMPONENTS OF A STRATIS VOLUME	160
32.3. BLOCK DEVICES USABLE WITH STRATIS	161
Supported devices	161
Unsupported devices	162
32.4. INSTALLING STRATIS	162
32.5. CREATING AN UNENCRYPTED STRATIS POOL	162
32.6. CREATING AN ENCRYPTED STRATIS POOL	163
32.7. SETTING OVERPROVISIONING MODE IN STRATIS FILESYSTEM	164
32.8. BINDING A STRATIS POOL TO NBDE	165
32.9. BINDING A STRATIS POOL TO TPM	166
32.10. UNLOCKING AN ENCRYPTED STRATIS POOL WITH KERNEL KEYRING	167
32.11. UNBINDING A STRATIS POOL FROM SUPPLEMENTARY ENCRYPTION	167
32.12. STARTING AND STOPPING STRATIS POOL	168
32.13. CREATING A STRATIS FILE SYSTEM	169
32.14. MOUNTING A STRATIS FILE SYSTEM	170
32.15. PERSISTENTLY MOUNTING A STRATIS FILE SYSTEM	170
32.16. SETTING UP NON-ROOT STRATIS FILESYSTEMS IN /ETC/FSTAB USING A SYSTEMD SERVICE	171
 CHAPTER 33. EXTENDING A STRATIS VOLUME WITH ADDITIONAL BLOCK DEVICES	 172
33.1. COMPONENTS OF A STRATIS VOLUME	172
33.2. ADDING BLOCK DEVICES TO A STRATIS POOL	173
33.3. ADDITIONAL RESOURCES	173
 CHAPTER 34. MONITORING STRATIS FILE SYSTEMS	 174
34.1. STRATIS SIZES REPORTED BY DIFFERENT UTILITIES	174
34.2. DISPLAYING INFORMATION ABOUT STRATIS VOLUMES	174
34.3. ADDITIONAL RESOURCES	175
 CHAPTER 35. USING SNAPSHOTS ON STRATIS FILE SYSTEMS	 176
35.1. CHARACTERISTICS OF STRATIS SNAPSHOTS	176
35.2. CREATING A STRATIS SNAPSHOT	176
35.3. ACCESSING THE CONTENT OF A STRATIS SNAPSHOT	177
35.4. REVERTING A STRATIS FILE SYSTEM TO A PREVIOUS SNAPSHOT	177
35.5. REMOVING A STRATIS SNAPSHOT	178
35.6. ADDITIONAL RESOURCES	178
 CHAPTER 36. REMOVING STRATIS FILE SYSTEMS	 179
36.1. COMPONENTS OF A STRATIS VOLUME	179
36.2. REMOVING A STRATIS FILE SYSTEM	180
36.3. REMOVING A STRATIS POOL	180
36.4. ADDITIONAL RESOURCES	181
 CHAPTER 37. GETTING STARTED WITH AN EXT4 FILE SYSTEM	 182
37.1. FEATURES OF AN EXT4 FILE SYSTEM	182
37.2. CREATING AN EXT4 FILE SYSTEM	182
37.3. MOUNTING AN EXT4 FILE SYSTEM	184
37.4. RESIZING AN EXT4 FILE SYSTEM	184
37.5. COMPARISON OF TOOLS USED WITH EXT4 AND XFS	185

MAKING OPEN SOURCE MORE INCLUSIVE

Red Hat is committed to replacing problematic language in our code, documentation, and web properties. We are beginning with these four terms: master, slave, blacklist, and whitelist. Because of the enormity of this endeavor, these changes will be implemented gradually over several upcoming releases. For more details, see [our CTO Chris Wright's message](#).

PROVIDING FEEDBACK ON RED HAT DOCUMENTATION

We appreciate your feedback on our documentation. Let us know how we can improve it.

Submitting comments on specific passages

1. View the documentation in the **Multi-page HTML** format and ensure that you see the **Feedback** button in the upper right corner after the page fully loads.
2. Use your cursor to highlight the part of the text that you want to comment on.
3. Click the **Add Feedback** button that appears near the highlighted text.
4. Add your feedback and click **Submit**.

Submitting feedback through Jira (account required)

1. Log in to the [Jira](#) website.
2. Click **Create** in the top navigation bar
3. Enter a descriptive title in the **Summary** field.
4. Enter your suggestion for improvement in the **Description** field. Include links to the relevant parts of the documentation.
5. Click **Create** at the bottom of the dialogue.

CHAPTER 1. OVERVIEW OF AVAILABLE FILE SYSTEMS

Choosing the file system that is appropriate for your application is an important decision due to the large number of options available and the trade-offs involved.

The following sections describe the file systems that Red Hat Enterprise Linux 9 includes by default, and recommendations on the most suitable file system for your application.

1.1. TYPES OF FILE SYSTEMS

Red Hat Enterprise Linux 9 supports a variety of file systems (FS). Different types of file systems solve different kinds of problems, and their usage is application specific. At the most general level, available file systems can be grouped into the following major types:

Table 1.1. Types of file systems and their use cases

Type	File system	Attributes and use cases
Disk or local FS	XFS	XFS is the default file system in RHEL. Red Hat recommends deploying XFS as your local file system unless there are specific reasons to do otherwise: for example, compatibility or corner cases around performance.
	ext4	ext4 has the benefit of familiarity in Linux, having evolved from the older ext2 and ext3 file systems. In many cases, it rivals XFS on performance. Support limits for ext4 filesystem and file sizes are lower than those on XFS.
Network or client-and-server FS	NFS	Use NFS to share files between multiple systems on the same network.
	SMB	Use SMB for file sharing with Microsoft Windows systems.
Shared storage or shared disk FS	GFS2	GFS2 provides shared write access to members of a compute cluster. The emphasis is on stability and reliability, with the functional experience of a local file system as possible. SAS Grid, Tibco MQ, IBM Websphere MQ, and Red Hat Active MQ have been deployed successfully on GFS2.
Volume-managing FS	Stratis (Technology Preview)	Stratis is a volume manager built on a combination of XFS and LVM. The purpose of Stratis is to emulate capabilities offered by volume-managing file systems like Btrfs and ZFS. It is possible to build this stack manually, but Stratis reduces configuration complexity, implements best practices, and consolidates error information.

1.2. LOCAL FILE SYSTEMS

Local file systems are file systems that run on a single, local server and are directly attached to storage.

For example, a local file system is the only choice for internal SATA or SAS disks, and is used when your server has internal hardware RAID controllers with local drives. Local file systems are also the most common file systems used on SAN attached storage when the device exported on the SAN is not shared.

All local file systems are POSIX-compliant and are fully compatible with all supported Red Hat Enterprise Linux releases. POSIX-compliant file systems provide support for a well-defined set of system calls, such as **read()**, **write()**, and **seek()**.

From the application programmer's point of view, there are relatively few differences between local file systems. The most notable differences from a user's perspective are related to scalability and performance. When considering a file system choice, consider how large the file system needs to be, what unique features it should have, and how it performs under your workload.

Available local file systems

- XFS
- ext4

1.3. THE XFS FILE SYSTEM

XFS is a highly scalable, high-performance, robust, and mature 64-bit journaling file system that supports very large files and file systems on a single host. It is the default file system in Red Hat Enterprise Linux 9. XFS was originally developed in the early 1990s by SGI and has a long history of running on extremely large servers and storage arrays.

The features of XFS include:

Reliability

- Metadata journaling, which ensures file system integrity after a system crash by keeping a record of file system operations that can be replayed when the system is restarted and the file system remounted
- Extensive run-time metadata consistency checking
- Scalable and fast repair utilities
- Quota journaling. This avoids the need for lengthy quota consistency checks after a crash.

Scalability and performance

- Supported file system size up to 1024 TiB
- Ability to support a large number of concurrent operations
- B-tree indexing for scalability of free space management
- Sophisticated metadata read-ahead algorithms
- Optimizations for streaming video workloads

Allocation schemes

- Extent-based allocation
- Stripe-aware allocation policies
- Delayed allocation
- Space pre-allocation
- Dynamically allocated inodes

Other features

- Replink-based file copies
- Tightly integrated backup and restore utilities
- Online defragmentation
- Online file system growing
- Comprehensive diagnostics capabilities
- Extended attributes (**xattr**). This allows the system to associate several additional name/value pairs per file.
- Project or directory quotas. This allows quota restrictions over a directory tree.
- Subsecond timestamps

Performance characteristics

XFS has a high performance on large systems with enterprise workloads. A large system is one with a relatively high number of CPUs, multiple HBAs, and connections to external disk arrays. XFS also performs well on smaller systems that have a multi-threaded, parallel I/O workload.

XFS has a relatively low performance for single threaded, metadata-intensive workloads: for example, a workload that creates or deletes large numbers of small files in a single thread.

1.4. THE EXT4 FILE SYSTEM

The ext4 file system is the fourth generation of the ext file system family. It was the default file system in Red Hat Enterprise Linux 6.

The ext4 driver can read and write to ext2 and ext3 file systems, but the ext4 file system format is not compatible with ext2 and ext3 drivers.

ext4 adds several new and improved features, such as:

- Supported file system size up to 50 TiB
- Extent-based metadata
- Delayed allocation

- Journal checksumming
- Large storage support

The extent-based metadata and the delayed allocation features provide a more compact and efficient way to track utilized space in a file system. These features improve file system performance and reduce the space consumed by metadata. Delayed allocation allows the file system to postpone selection of the permanent location for newly written user data until the data is flushed to disk. This enables higher performance since it can allow for larger, more contiguous allocations, allowing the file system to make decisions with much better information.

File system repair time using the **fsck** utility in ext4 is much faster than in ext2 and ext3. Some file system repairs have demonstrated up to a six-fold increase in performance.

1.5. COMPARISON OF XFS AND EXT4

XFS is the default file system in RHEL. This section compares the usage and features of XFS and ext4.

Metadata error behavior

In ext4, you can configure the behavior when the file system encounters metadata errors. The default behavior is to simply continue the operation. When XFS encounters an unrecoverable metadata error, it shuts down the file system and returns the **EFSCORRUPTED** error.

Quotas

In ext4, you can enable quotas when creating the file system or later on an existing file system. You can then configure the quota enforcement using a mount option.

XFS quotas are not a remountable option. You must activate quotas on the initial mount.

Running the **quotacheck** command on an XFS file system has no effect. The first time you turn on quota accounting, XFS checks quotas automatically.

File system resize

XFS has no utility to reduce the size of a file system. You can only increase the size of an XFS file system. In comparison, ext4 supports both extending and reducing the size of a file system.

Inode numbers

The ext4 file system does not support more than 2^{32} inodes.

XFS dynamically allocates inodes. An XFS file system cannot run out of inodes as long as there is free space on the file system.

Certain applications cannot properly handle inode numbers larger than 2^{32} on an XFS file system. These applications might cause the failure of 32-bit stat calls with the **Eoverflow** return value.

Inode number exceed 2^{32} under the following conditions:

- The file system is larger than 1 TiB with 256-byte inodes.
- The file system is larger than 2 TiB with 512-byte inodes.

If your application fails with large inode numbers, mount the XFS file system with the **-o inode32** option to enforce inode numbers below 2^{32} . Note that using **inode32** does not affect inodes that are already allocated with 64-bit numbers.



IMPORTANT

Do *not* use the **inode32** option unless a specific environment requires it. The **inode32** option changes allocation behavior. As a consequence, the **ENOSPC** error might occur if no space is available to allocate inodes in the lower disk blocks.

1.6. CHOOSING A LOCAL FILE SYSTEM

To choose a file system that meets your application requirements, you need to understand the target system on which you are going to deploy the file system. You can use the following questions to inform your decision:

- Do you have a large server?
- Do you have large storage requirements or have a local, slow SATA drive?
- What kind of I/O workload do you expect your application to present?
- What are your throughput and latency requirements?
- How stable is your server and storage hardware?
- What is the typical size of your files and data set?
- If the system fails, how much downtime can you suffer?

If both your server and your storage device are large, XFS is the best choice. Even with smaller storage arrays, XFS performs very well when the average file sizes are large (for example, hundreds of megabytes in size).

If your existing workload has performed well with ext4, staying with ext4 should provide you and your applications with a very familiar environment.

The ext4 file system tends to perform better on systems that have limited I/O capability. It performs better on limited bandwidth (less than 200MB/s) and up to around 1000 IOPS capability. For anything with higher capability, XFS tends to be faster.

XFS consumes about twice the CPU-per-metadata operation compared to ext4, so if you have a CPU-bound workload with little concurrency, then ext4 will be faster. In general, ext4 is better if an application uses a single read/write thread and small files, while XFS shines when an application uses multiple read/write threads and bigger files.

You cannot shrink an XFS file system. If you need to be able to shrink the file system, consider using ext4, which supports offline shrinking.

In general, Red Hat recommends that you use XFS unless you have a specific use case for ext4. You should also measure the performance of your specific application on your target server and storage system to make sure that you choose the appropriate type of file system.

Table 1.2. Summary of local file system recommendations

Scenario	Recommended file system
No special use case	XFS

Scenario	Recommended file system
Large server	XFS
Large storage devices	XFS
Large files	XFS
Multi-threaded I/O	XFS
Single-threaded I/O	ext4
Limited I/O capability (under 1000 IOPS)	ext4
Limited bandwidth (under 200MB/s)	ext4
CPU-bound workload	ext4
Support for offline shrinking	ext4

1.7. NETWORK FILE SYSTEMS

Network file systems, also referred to as client/server file systems, enable client systems to access files that are stored on a shared server. This makes it possible for multiple users on multiple systems to share files and storage resources.

Such file systems are built from one or more servers that export a set of file systems to one or more clients. The client nodes do not have access to the underlying block storage, but rather interact with the storage using a protocol that allows for better access control.

Available network file systems

- The most common client/server file system for RHEL customers is the NFS file system. RHEL provides both an NFS server component to export a local file system over the network and an NFS client to import these file systems.
- RHEL also includes a CIFS client that supports the popular Microsoft SMB file servers for Windows interoperability. The userspace Samba server provides Windows clients with a Microsoft SMB service from a RHEL server.

1.8. SHARED STORAGE FILE SYSTEMS

Shared storage file systems, sometimes referred to as cluster file systems, give each server in the cluster direct access to a shared block device over a local storage area network (SAN).

Comparison with network file systems

Like client/server file systems, shared storage file systems work on a set of servers that are all members of a cluster. Unlike NFS, however, no single server provides access to data or metadata to other members: each member of the cluster has direct access to the same storage device (the

shared storage), and all cluster member nodes access the same set of files.

Concurrency

Cache coherency is key in a clustered file system to ensure data consistency and integrity. There must be a single version of all files in a cluster visible to all nodes within a cluster. The file system must prevent members of the cluster from updating the same storage block at the same time and causing data corruption. In order to do that, shared storage file systems use a cluster wide-locking mechanism to arbitrate access to the storage as a concurrency control mechanism. For example, before creating a new file or writing to a file that is opened on multiple servers, the file system component on the server must obtain the correct lock.

The requirement of cluster file systems is to provide a highly available service like an Apache web server. Any member of the cluster will see a fully coherent view of the data stored in their shared disk file system, and all updates will be arbitrated correctly by the locking mechanisms.

Performance characteristics

Shared disk file systems do not always perform as well as local file systems running on the same system due to the computational cost of the locking overhead. Shared disk file systems perform well with workloads where each node writes almost exclusively to a particular set of files that are not shared with other nodes or where a set of files is shared in an almost exclusively read-only manner across a set of nodes. This results in a minimum of cross-node cache invalidation and can maximize performance.

Setting up a shared disk file system is complex, and tuning an application to perform well on a shared disk file system can be challenging.

Available shared storage file systems

- Red Hat Enterprise Linux provides the GFS2 file system. GFS2 comes tightly integrated with the Red Hat Enterprise Linux High Availability Add-On and the Resilient Storage Add-On.

Red Hat Enterprise Linux supports GFS2 on clusters that range in size from 2 to 16 nodes.

1.9. CHOOSING BETWEEN NETWORK AND SHARED STORAGE FILE SYSTEMS

When choosing between network and shared storage file systems, consider the following points:

- NFS-based network file systems are an extremely common and popular choice for environments that provide NFS servers.
- Network file systems can be deployed using very high-performance networking technologies like Infiniband or 10 Gigabit Ethernet. This means that you should not turn to shared storage file systems just to get raw bandwidth to your storage. If the speed of access is of prime importance, then use NFS to export a local file system like XFS.
- Shared storage file systems are not easy to set up or to maintain, so you should deploy them only when you cannot provide your required availability with either local or network file systems.
- A shared storage file system in a clustered environment helps reduce downtime by eliminating the steps needed for unmounting and mounting that need to be done during a typical fail-over scenario involving the relocation of a high-availability service.

Red Hat recommends that you use network file systems unless you have a specific use case for shared storage file systems. Use shared storage file systems primarily for deployments that need to provide high-availability services with minimum downtime and have stringent service-level requirements.

1.10. VOLUME-MANAGING FILE SYSTEMS

Volume-managing file systems integrate the entire storage stack for the purposes of simplicity and in-stack optimization.

Available volume-managing file systems

- Red Hat Enterprise Linux 9 provides the Stratis volume manager as a Technology Preview. Stratis uses XFS for the file system layer and integrates it with LVM, Device Mapper, and other components.

Stratis was first released in Red Hat Enterprise Linux 8.0. It is conceived to fill the gap created when Red Hat deprecated Btrfs. Stratis 1.0 is an intuitive, command line-based volume manager that can perform significant storage management operations while hiding the complexity from the user:

- Volume management
- Pool creation
- Thin storage pools
- Snapshots
- Automated read cache

Stratis offers powerful features, but currently lacks certain capabilities of other offerings that it might be compared to, such as Btrfs or ZFS. Most notably, it does not support CRCs with self healing.

CHAPTER 2. MANAGING LOCAL STORAGE USING RHEL SYSTEM ROLES

To manage LVM and local file systems (FS) using Ansible, you can use the **storage** role, which is one of the RHEL System Roles available in RHEL 9.

Using the **storage** role enables you to automate administration of file systems on disks and logical volumes on multiple machines and across all versions of RHEL starting with RHEL 7.7.

For more information about RHEL System Roles and how to apply them, see [Introduction to RHEL System Roles](#).

2.1. INTRODUCTION TO THE **STORAGE** RHEL SYSTEM ROLE

The **storage** role can manage:

- File systems on disks which have not been partitioned
- Complete LVM volume groups including their logical volumes and file systems
- MD RAID volumes and their file systems

With the **storage** role, you can perform the following tasks:

- Create a file system
- Remove a file system
- Mount a file system
- Unmount a file system
- Create LVM volume groups
- Remove LVM volume groups
- Create logical volumes
- Remove logical volumes
- Create RAID volumes
- Remove RAID volumes
- Create LVM volume groups with RAID
- Remove LVM volume groups with RAID
- Create encrypted LVM volume groups
- Create LVM logical volumes with RAID

2.2. PARAMETERS THAT IDENTIFY A STORAGE DEVICE IN THE **STORAGE** RHEL SYSTEM ROLE

Your **storage** role configuration affects only the file systems, volumes, and pools that you list in the following variables.

storage_volumes

List of file systems on all unpartitioned disks to be managed.

storage_volumes can also include **raid** volumes.

Partitions are currently unsupported.

storage_pools

List of pools to be managed.

Currently the only supported pool type is LVM. With LVM, pools represent volume groups (VGs).

Under each pool there is a list of volumes to be managed by the role. With LVM, each volume corresponds to a logical volume (LV) with a file system.

2.3. EXAMPLE ANSIBLE PLAYBOOK TO CREATE AN XFS FILE SYSTEM ON A BLOCK DEVICE

The example Ansible playbook applies the **storage** role to create an XFS file system on a block device using the default parameters.



WARNING

The **storage** role can create a file system only on an unpartitioned, whole disk or a logical volume (LV). It cannot create the file system on a partition.

Example 2.1. A playbook that creates XFS on /dev/sdb

```
---
- hosts: all
  vars:
    storage_volumes:
      - name: barefs
        type: disk
        disks:
          - sdb
        fs_type: xfs
  roles:
    - rhel-system-roles.storage
```

- The volume name (**barefs** in the example) is currently arbitrary. The **storage** role identifies the volume by the disk device listed under the **disks:** attribute.
- You can omit the **fs_type: xfs** line because XFS is the default file system in RHEL 9.

- To create the file system on an LV, provide the LVM setup under the **disks:** attribute, including the enclosing volume group. For details, see [Example Ansible playbook to manage logical volumes](#).
Do not provide the path to the LV device.

Additional resources

- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file.

2.4. EXAMPLE ANSIBLE PLAYBOOK TO PERSISTENTLY MOUNT A FILE SYSTEM

The example Ansible applies the **storage** role to immediately and persistently mount an XFS file system.

Example 2.2. A playbook that mounts a file system on `/dev/sdb` to `/mnt/data`

```
---
- hosts: all
  vars:
    storage_volumes:
      - name: barefs
        type: disk
        disks:
          - sdb
        fs_type: xfs
        mount_point: /mnt/data
  roles:
    - rhel-system-roles.storage
```

- This playbook adds the file system to the `/etc/fstab` file, and mounts the file system immediately.
- If the file system on the `/dev/sdb` device or the mount point directory do not exist, the playbook creates them.

Additional resources

- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file.

2.5. EXAMPLE ANSIBLE PLAYBOOK TO MANAGE LOGICAL VOLUMES

The example Ansible playbook applies the **storage** role to create an LVM logical volume in a volume group.

Example 2.3. A playbook that creates a mylv logical volume in the myvg volume group

```
- hosts: all
  vars:
    storage_pools:
      - name: myvg
        disks:
```



```

- sda
- sdb
- sdc
volumes:
- name: mylv
  size: 2G
  fs_type: ext4
  mount_point: /mnt/data
roles:
- rhel-system-roles.storage

```

- The **myvg** volume group consists of the following disks:
 - **/dev/sda**
 - **/dev/sdb**
 - **/dev/sdc**
- If the **myvg** volume group already exists, the playbook adds the logical volume to the volume group.
- If the **myvg** volume group does not exist, the playbook creates it.
- The playbook creates an Ext4 file system on the **mylv** logical volume, and persistently mounts the file system at **/mnt**.

Additional resources

- The **/usr/share/ansible/roles/rhel-system-roles.storage/README.md** file.

2.6. EXAMPLE ANSIBLE PLAYBOOK TO ENABLE ONLINE BLOCK DISCARD

The example Ansible playbook applies the **storage** role to mount an XFS file system with online block discard enabled.

Example 2.4. A playbook that enables online block discard on /mnt/data/

```

---
- hosts: all
  vars:
    storage_volumes:
      - name: barefs
        type: disk
        disks:
          - sdb
        fs_type: xfs
        mount_point: /mnt/data
        mount_options: discard
  roles:
    - rhel-system-roles.storage

```

Additional resources

- [Example Ansible playbook to persistently mount a file system](#)
- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file.

2.7. EXAMPLE ANSIBLE PLAYBOOK TO CREATE AND MOUNT AN EXT4 FILE SYSTEM

The example Ansible playbook applies the **storage** role to create and mount an Ext4 file system.

Example 2.5. A playbook that creates Ext4 on `/dev/sdb` and mounts it at `/mnt/data`

```
---
- hosts: all
  vars:
    storage_volumes:
      - name: barefs
        type: disk
        disks:
          - sdb
        fs_type: ext4
        fs_label: label-name
        mount_point: /mnt/data
  roles:
    - rhel-system-roles.storage
```

- The playbook creates the file system on the `/dev/sdb` disk.
- The playbook persistently mounts the file system at the `/mnt/data` directory.
- The label of the file system is ***label-name***.

Additional resources

- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file.

2.8. EXAMPLE ANSIBLE PLAYBOOK TO CREATE AND MOUNT AN EXT3 FILE SYSTEM

The example Ansible playbook applies the **storage** role to create and mount an Ext3 file system.

Example 2.6. A playbook that creates Ext3 on `/dev/sdb` and mounts it at `/mnt/data`

```
---
- hosts: all
  vars:
    storage_volumes:
      - name: barefs
        type: disk
        disks:
          - sdb
```

```

fs_type: ext3
fs_label: label-name
mount_point: /mnt/data
roles:
  - rhel-system-roles.storage

```

- The playbook creates the file system on the **/dev/sdb** disk.
- The playbook persistently mounts the file system at the **/mnt/data** directory.
- The label of the file system is ***label-name***.

Additional resources

- The **/usr/share/ansible/roles/rhel-system-roles.storage/README.md** file.

2.9. EXAMPLE ANSIBLE PLAYBOOK TO RESIZE AN EXISTING EXT4 OR EXT3 FILE SYSTEM USING THE STORAGE RHEL SYSTEM ROLE

The example Ansible playbook applies the **storage** role to resize an existing Ext4 or Ext3 file system on a block device.

Example 2.7. A playbook that set up a single volume on a disk

```

---
- name: Create a disk device mounted on /opt/barefs
  hosts: all
  vars:
    storage_volumes:
      - name: barefs
        type: disk
        disks:
          - /dev/sdb
        size: 12 GiB
        fs_type: ext4
        mount_point: /opt/barefs
  roles:
    - rhel-system-roles.storage

```

- If the volume in the previous example already exists, to resize the volume, you need to run the same playbook, just with a different value for the parameter **size**. For example:

Example 2.8. A playbook that resizes ext4 on /dev/sdb

```

---
- name: Create a disk device mounted on /opt/barefs
  hosts: all
  vars:
    storage_volumes:
      - name: barefs
        type: disk

```

```

disks:
  - /dev/sdb
size: 10 GiB
fs_type: ext4
mount_point: /opt/barefs
roles:
  - rhel-system-roles.storage

```

- The volume name (barefs in the example) is currently arbitrary. The Storage role identifies the volume by the disk device listed under the disks: attribute.



NOTE

Using the **Resizing** action in other file systems can destroy the data on the device you are working on.

Additional resources

- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file.

2.10. EXAMPLE ANSIBLE PLAYBOOK TO RESIZE AN EXISTING FILE SYSTEM ON LVM USING THE STORAGE RHEL SYSTEM ROLE

The example Ansible playbook applies the **storage** RHEL System Role to resize an LVM logical volume with a file system.



WARNING

Using the **Resizing** action in other file systems can destroy the data on the device you are working on.

Example 2.9. A playbook that resizes existing mylv1 and mylv2 logical volumes in the myvg volume group

```

---

- hosts: all
  vars:
    storage_pools:
      - name: myvg
        disks:
          - /dev/sda
          - /dev/sdb
          - /dev/sdc
        volumes:
          - name: mylv1
            size: 10 GiB

```

```

    fs_type: ext4
    mount_point: /opt/mount1
  - name: mylv2
    size: 50 GiB
    fs_type: ext4
    mount_point: /opt/mount2

```

```

- name: Create LVM pool over three disks
  include_role:
    name: rhel-system-roles.storage

```

- This playbook resizes the following existing file systems:
 - The Ext4 file system on the **mylv1** volume, which is mounted at **/opt/mount1**, resizes to 10 GiB.
 - The Ext4 file system on the **mylv2** volume, which is mounted at **/opt/mount2**, resizes to 50 GiB.

Additional resources

- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file.

2.11. EXAMPLE ANSIBLE PLAYBOOK TO CREATE A SWAP VOLUME USING THE STORAGE RHEL SYSTEM ROLE

This section provides an example Ansible playbook. This playbook applies the **storage** role to create a swap volume, if it does not exist, or to modify the swap volume, if it already exist, on a block device using the default parameters.

Example 2.10. A playbook that creates or modify an existing XFS on /dev/sdb

```

---
- name: Create a disk device with swap
  hosts: all
  vars:
    storage_volumes:
      - name: swap_fs
        type: disk
        disks:
          - /dev/sdb
    size: 15 GiB
    fs_type: swap
  roles:
    - rhel-system-roles.storage

```

- The volume name (**swap_fs** in the example) is currently arbitrary. The **storage** role identifies the volume by the disk device listed under the **disks:** attribute.

Additional resources

- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file.

2.12. CONFIGURING A RAID VOLUME USING THE STORAGE SYSTEM ROLE

With the **storage** System Role, you can configure a RAID volume on RHEL using Red Hat Ansible Automation Platform and Ansible-Core. Create an Ansible playbook with the parameters to configure a RAID volume to suit your requirements.

Prerequisites

- The Ansible Core package is installed on the control machine.
- You have the **rhel-system-roles** package installed on the system from which you want to run the playbook.
- You have an inventory file detailing the systems on which you want to deploy a RAID volume using the **storage** System Role.

Procedure

1. Create a new *playbook.yml* file with the following content:

```
---
- name: Configure the storage
  hosts: managed-node-01.example.com
  tasks:
  - name: Create a RAID on sdd, sde, sdf, and sdg
    include_role:
      name: rhel-system-roles.storage
  vars:
    storage_safe_mode: false
    storage_volumes:
      - name: data
        type: raid
        disks: [sdd, sde, sdf, sdg]
        raid_level: raid0
        raid_chunk_size: 32 KiB
        mount_point: /mnt/data
        state: present
```



WARNING

Device names might change in certain circumstances, for example, when you add a new disk to a system. Therefore, to prevent data loss, do not use specific disk names in the playbook.

2. Optional: Verify the playbook syntax:

```
# ansible-playbook --syntax-check playbook.yml
```

3. Run the playbook:

```
# ansible-playbook -i inventory.file /path/to/file/playbook.yml
```

Additional resources

- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file
- [Preparing a control node and managed nodes to use RHEL System Roles](#)

2.13. CONFIGURING AN LVM POOL WITH RAID USING THE `storage` RHEL SYSTEM ROLE

With the **storage** System Role, you can configure an LVM pool with RAID on RHEL using Red Hat Ansible Automation Platform. You can set up an Ansible playbook with the available parameters to configure an LVM pool with RAID.

Prerequisites

- The Ansible Core package is installed on the control machine.
- You have the **rhel-system-roles** package installed on the system from which you want to run the playbook.
- You have an inventory file detailing the systems on which you want to configure an LVM pool with RAID using the **storage** System Role.

Procedure

1. Create a new **playbook.yml** file with the following content:

```
- hosts: all
  vars:
    storage_safe_mode: false
    storage_pools:
      - name: my_pool
        type: lvm
        disks: [sdh, sdi]
        raid_level: raid1
        volumes:
          - name: my_pool
            size: "1 GiB"
            mount_point: "/mnt/app/shared"
            fs_type: xfs
            state: present
  roles:
    - name: rhel-system-roles.storage
```



NOTE

To create an LVM pool with RAID, you must specify the RAID type using the **raid_level** parameter.

- Optional. Verify playbook syntax.

```
# ansible-playbook --syntax-check playbook.yml
```

- Run the playbook on your inventory file:

```
# ansible-playbook -i inventory.file /path/to/file/playbook.yml
```

Additional resources

- The `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file.

2.14. EXAMPLE ANSIBLE PLAYBOOK TO COMPRESS AND DEDUPLICATE A VDO VOLUME ON LVM USING THE STORAGE RHEL SYSTEM ROLE

The example Ansible playbook applies the **storage** RHEL System Role to enable compression and deduplication of Logical Volumes (LVM) using Virtual Data Optimizer (VDO).

Example 2.11. A playbook that creates `mylv1` LVM VDO volume in the `myvg` volume group

```
---
- name: Create LVM VDO volume under volume group 'myvg'
  hosts: all
  roles:
    - rhel-system-roles.storage
  vars:
    storage_pools:
      - name: myvg
        disks:
          - /dev/sdb
        volumes:
          - name: mylv1
            compression: true
            deduplication: true
            vdo_pool_size: 10 GiB
            size: 30 GiB
            mount_point: /mnt/app/shared
```

In this example, the **compression** and **deduplication** pools are set to true, which specifies that the VDO is used. The following describes the usage of these parameters:

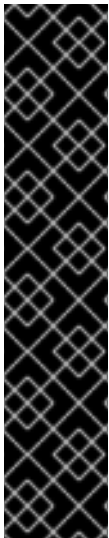
- The **deduplication** is used to deduplicate the duplicated data stored on the storage volume.
- The compression is used to compress the data stored on the storage volume, which results in more storage capacity.
- The `vdo_pool_size` specifies the actual size the volume takes on the device. The virtual size of VDO volume is set by the **size** parameter. NOTE: Because of the Storage role use of LVM VDO, only one volume per pool can use the compression and deduplication.

2.15. CREATING A LUKS2 ENCRYPTED VOLUME USING THE STORAGE RHEL SYSTEM ROLE

You can use the **storage** role to create and configure a volume encrypted with LUKS by running an Ansible playbook.

Prerequisites

- Access and permissions to one or more managed nodes, which are systems you want to configure with the **crypto_policies** System Role.
- An inventory file, which lists the managed nodes.
- Access and permissions to a control node, which is a system from which Red Hat Ansible Core configures other systems. On the control node, the **ansible-core** and **rhel-system-roles** packages are installed.



IMPORTANT

RHEL 8.0–8.5 provided access to a separate Ansible repository that contains Ansible Engine 2.9 for automation based on Ansible. Ansible Engine contains command-line utilities such as **ansible**, **ansible-playbook**, connectors such as **docker** and **podman**, and many plugins and modules. For information about how to obtain and install Ansible Engine, see the [How to download and install Red Hat Ansible Engine](#) Knowledgebase article.

RHEL 8.6 and 9.0 have introduced Ansible Core (provided as the **ansible-core** package), which contains the Ansible command-line utilities, commands, and a small set of built-in Ansible plugins. RHEL provides this package through the AppStream repository, and it has a limited scope of support. For more information, see the [Scope of support for the Ansible Core package included in the RHEL 9 and RHEL 8.6 and later AppStream repositories](#) Knowledgebase article.

Procedure

1. Create a new **playbook.yml** file with the following content:

```
- hosts: all
  vars:
    storage_volumes:
      - name: barefs
        type: disk
        disks:
          - sdb
        fs_type: xfs
        fs_label: label-name
        mount_point: /mnt/data
        encryption: true
        encryption_password: your-password
  roles:
    - rhel-system-roles.storage
```

You can also add the other encryption parameters such as **encryption_key**, **encryption_cipher**, **encryption_key_size**, and **encryption_luks** version in the **playbook.yml** file.

2. Optional: Verify playbook syntax:

```
# ansible-playbook --syntax-check playbook.yml
```

3. Run the playbook on your inventory file:

```
# ansible-playbook -i inventory.file /path/to/file/playbook.yml
```

Verification

1. View the encryption status:

```
# cryptsetup status sdb

/dev/mapper/sdb is active and is in use.
type: LUKS2
cipher: aes-xts-plain64
keysize: 512 bits
key location: keyring
device: /dev/sdb
[...]
```

2. Verify the created LUKS encrypted volume:

```
# cryptsetup luksDump /dev/sdb

Version:      2
Epoch:       6
Metadata area: 16384 [bytes]
Keyslots area: 33521664 [bytes]
UUID:        a4c6be82-7347-4a91-a8ad-9479b72c9426
Label:       (no label)
Subsystem:   (no subsystem)
Flags:       allow-discards

Data segments:
 0: crypt
  offset: 33554432 [bytes]
  length: (whole device)
  cipher: aes-xts-plain64
  sector: 4096 [bytes]
[...]
```

3. View the **cryptsetup** parameters in the **playbook.yml** file, which the **storage** role supports:

```
# cat ~/playbook.yml

- hosts: all
  vars:
    storage_volumes:
      - name: foo
        type: disk
        disks:
          - nvme0n1
```

```

    fs_type: xfs
    fs_label: label-name
    mount_point: /mnt/data
    encryption: true
    #encryption_password: passwdpasswd
    encryption_key: /home/passwd_key
    encryption_cipher: aes-xts-plain64
    encryption_key_size: 512
    encryption_luks_version: luks2

roles:
  - rhel-system-roles.storage

```

Additional resources

- [Encrypting block devices using LUKS](#)
- `/usr/share/ansible/roles/rhel-system-roles.storage/README.md` file

2.16. EXAMPLE ANSIBLE PLAYBOOK TO EXPRESS POOL VOLUME SIZES AS PERCENTAGE USING THE STORAGE RHEL SYSTEM ROLE

The example Ansible playbook applies the **storage** System Role to enable you to express Logical Manager Volumes (LVM) volume sizes as a percentage of the pool's total size.

Example 2.12. A playbook that express volume sizes as a percentage of the pool's total size

```

---
- name: Express volume sizes as a percentage of the pool's total size
  hosts: all
  roles:
    - rhel-system-roles.storage
  vars:
    storage_pools:
      - name: myvg
        disks:
          - /dev/sdb
        volumes:
          - name: data
            size: 60%
            mount_point: /opt/mount/data
          - name: web
            size: 30%
            mount_point: /opt/mount/web
          - name: cache
            size: 10%
            mount_point: /opt/cache/mount

```

This example specifies the size of LVM volumes as a percentage of the pool size, for example: "60%". Additionally, you can also specify the size of LVM volumes as a percentage of the pool size in a human-readable size of the file system, for example, "10g" or "50 GiB".

2.17. ADDITIONAL RESOURCES

- [/usr/share/doc/rhel-system-roles/storage/](#)
- [/usr/share/ansible/roles/rhel-system-roles.storage/](#)

CHAPTER 3. MOUNTING NFS SHARES

As a system administrator, you can mount remote NFS shares on your system to access shared data.

3.1. INTRODUCTION TO NFS

This section explains the basic concepts of the NFS service.

A Network File System (NFS) allows remote hosts to mount file systems over a network and interact with those file systems as though they are mounted locally. This enables you to consolidate resources onto centralized servers on the network.

The NFS server refers to the **/etc/exports** configuration file to determine whether the client is allowed to access any exported file systems. Once verified, all file and directory operations are available to the user.

3.2. SUPPORTED NFS VERSIONS

This section lists versions of NFS supported in Red Hat Enterprise Linux and their features.

Currently, Red Hat Enterprise Linux 9 supports the following major versions of NFS:

- NFS version 3 (NFSv3) supports safe asynchronous writes and is more robust at error handling than the previous NFSv2; it also supports 64-bit file sizes and offsets, allowing clients to access more than 2 GB of file data.
- NFS version 4 (NFSv4) works through firewalls and on the Internet, no longer requires an **rpcbind** service, supports Access Control Lists (ACLs), and utilizes stateful operations.

NFS version 2 (NFSv2) is no longer supported by Red Hat.

Default NFS version

The default NFS version in Red Hat Enterprise Linux 9 is 4.2. NFS clients attempt to mount using NFSv4.2 by default, and fall back to NFSv4.1 when the server does not support NFSv4.2. The mount later falls back to NFSv4.0 and then to NFSv3.

Features of minor NFS versions

Following are the features of NFSv4.2 in Red Hat Enterprise Linux 9:

Server-side copy

Enables the NFS client to efficiently copy data without wasting network resources using the **copy_file_range()** system call.

Sparse files

Enables files to have one or more *holes*, which are unallocated or uninitialized data blocks consisting only of zeroes. The **lseek()** operation in NFSv4.2 supports **seek_hole()** and **seek_data()**, which enables applications to map out the location of holes in the sparse file.

Space reservation

Permits storage servers to reserve free space, which prohibits servers to run out of space. NFSv4.2 supports the **allocate()** operation to reserve space, the **deallocate()** operation to unreserve space, and the **fallocate()** operation to preallocate or deallocate space in a file.

Labeled NFS

Enforces data access rights and enables SELinux labels between a client and a server for individual files on an NFS file system.

Layout enhancements

Provides the **layoutstats()** operation, which enables some Parallel NFS (pNFS) servers to collect better performance statistics.

Following are the features of NFSv4.1:

- Enhances performance and security of network, and also includes client-side support for pNFS.
- No longer requires a separate TCP connection for callbacks, which allows an NFS server to grant delegations even when it cannot contact the client: for example, when NAT or a firewall interferes.
- Provides exactly once semantics (except for reboot operations), preventing a previous issue whereby certain operations sometimes returned an inaccurate result if a reply was lost and the operation was sent twice.

3.3. SERVICES REQUIRED BY NFS

This section lists system services that are required for running an NFS server or mounting NFS shares. Red Hat Enterprise Linux starts these services automatically.

Red Hat Enterprise Linux uses a combination of kernel-level support and service processes to provide NFS file sharing. All NFS versions rely on Remote Procedure Calls (RPC) between clients and servers. To share or mount NFS file systems, the following services work together depending on which version of NFS is implemented:

nfsd

The NFS server kernel module that services requests for shared NFS file systems.

rpcbind

Accepts port reservations from local RPC services. These ports are then made available (or advertised) so the corresponding remote RPC services can access them. The **rpcbind** service responds to requests for RPC services and sets up connections to the requested RPC service. This is not used with NFSv4.

rpc.mountd

This process is used by an NFS server to process **MOUNT** requests from NFSv3 clients. It checks that the requested NFS share is currently exported by the NFS server, and that the client is allowed to access it. If the mount request is allowed, the **nfs-mountd** service replies with a Success status and provides the File-Handle for this NFS share back to the NFS client.

rpc.nfsd

This process enables explicit NFS versions and protocols the server advertises to be defined. It works with the Linux kernel to meet the dynamic demands of NFS clients, such as providing server threads each time an NFS client connects. This process corresponds to the **nfs-server** service.

lockd

This is a kernel thread that runs on both clients and servers. It implements the Network Lock Manager (NLM) protocol, which enables NFSv3 clients to lock files on the server. It is started automatically whenever the NFS server is run and whenever an NFS file system is mounted.

rpc.statd

This process implements the Network Status Monitor (NSM) RPC protocol, which notifies NFS clients when an NFS server is restarted without being gracefully brought down. The **rpc-statd** service is started automatically by the **nfs-server** service, and does not require user configuration. This is not used with NFSv4.

rpc.rquotad

This process provides user quota information for remote users. The **rpc-rquotad** service, which is provided by the **quota-rpc** package, has to be started by user when the **nfs-server** is started.

rpc.idmapd

This process provides NFSv4 client and server upcalls, which map between on-the-wire NFSv4 names (strings in the form of **user@domain**) and local UIDs and GIDs. For **idmapd** to function with NFSv4, the **/etc/idmapd.conf** file must be configured. At a minimum, the **Domain** parameter should be specified, which defines the NFSv4 mapping domain. If the NFSv4 mapping domain is the same as the DNS domain name, this parameter can be skipped. The client and server must agree on the NFSv4 mapping domain for ID mapping to function properly.

Only the NFSv4 server uses **rpc.idmapd**, which is started by the **nfs-idmapd** service. The NFSv4 client uses the keyring-based **nfsidmap** utility, which is called by the kernel on-demand to perform ID mapping. If there is a problem with **nfsidmap**, the client falls back to using **rpc.idmapd**.

The RPC services with NFSv4

The mounting and locking protocols have been incorporated into the NFSv4 protocol. The server also listens on the well-known TCP port 2049. As such, NFSv4 does not need to interact with **rpcbind**, **lockd**, and **rpc-statd** services. The **nfs-mountd** service is still required on the NFS server to set up the exports, but is not involved in any over-the-wire operations.

Additional resources

- [Configuring an NFSv4 only server without **rpcbind**](#).

3.4. NFS HOST NAME FORMATS

This section describes different formats that you can use to specify a host when mounting or exporting an NFS share.

You can specify the host in the following formats:

Single machine

Either of the following:

- A fully-qualified domain name (that can be resolved by the server)
- Host name (that can be resolved by the server)
- An IP address.

IP networks

Either of the following formats is valid:

- **a.b.c.d/z**, where **a.b.c.d** is the network and **z** is the number of bits in the netmask; for example **192.168.0.0/24**.
- **a.b.c.d/netmask**, where **a.b.c.d** is the network and **netmask** is the netmask; for example, **192.168.100.8/255.255.255.0**.

Netgroups

The **@group-name** format, where **group-name** is the NIS netgroup name.

3.5. INSTALLING NFS

This procedure installs all packages necessary to mount or export NFS shares.

Procedure

- Install the **nfs-utils** package:

```
# dnf install nfs-utils
```

3.6. DISCOVERING NFS EXPORTS

This procedure discovers which file systems a given NFSv3 or NFSv4 server exports.

Procedure

- With any server that supports NFSv3, use the **showmount** utility:

```
$ showmount --exports my-server
```

```
Export list for my-server
```

```
/exports/foo
```

```
/exports/bar
```

- With any server that supports NFSv4, mount the root directory and look around:

```
# mount my-server:/ /mnt/
```

```
# ls /mnt/
```

```
exports
```

```
# ls /mnt/exports/
```

```
foo
```

```
bar
```

On servers that support both NFSv4 and NFSv3, both methods work and give the same results.

Additional resources

- **showmount(8)** man page

3.7. MOUNTING AN NFS SHARE WITH MOUNT

Mount an NFS share exported from a server by using the **mount** utility.



WARNING

You can experience conflicts in your NFSv4 **clientid** and their sudden expiration if your NFS clients have the same short hostname. To avoid any possible sudden expiration of your NFSv4 **clientid**, you must use either unique hostnames for NFS clients or configure identifier on each container, depending on what system you are using. For more information, see the [NFSv4 clientid was expired suddenly due to use same hostname on several NFS clients](#) Knowledgebase article.

Procedure

- To mount an NFS share, use the following command:

```
# mount -t nfs -o options host:/remote/export /local/directory
```

This command uses the following variables:

options

A comma-delimited list of mount options.

host

The host name, IP address, or fully qualified domain name of the server exporting the file system you want to mount.

/remote/export

The file system or directory being exported from the server, that is, the directory you want to mount.

/local/directory

The client location where */remote/export* is mounted.

Additional resources

- [Common NFS mount options.](#)
- [NFS host name formats .](#)
- [Mounting a file system with mount .](#)
- mount(8)** man page
- exports(5)** man page

3.8. COMMON NFS MOUNT OPTIONS

The following are the commonly used options when mounting NFS shares. You can use these options with manual **mount** commands, the **/etc/fstab** settings, and **autofs**.

lookupcache=mode

Specifies how the kernel should manage its cache of directory entries for a given mount point. Valid arguments for *mode* are **all**, **none**, or **positive**.

nfsvers=version

Specifies which version of the NFS protocol to use, where *version* is **3**, **4**, **4.0**, **4.1**, or **4.2**. This is useful for hosts that run multiple NFS servers, or to disable retrying a mount with lower versions. If no version is specified, NFS uses the highest version supported by the kernel and the **mount** utility. The option **vers** is identical to **nfsvers**, and is included in this release for compatibility reasons.

noacl

Turns off all ACL processing. This may be needed when interfacing with older versions of Red Hat Enterprise Linux, Red Hat Linux, or Solaris, because the most recent ACL technology is not compatible with older systems.

noLOCK

Disables file locking. This setting is sometimes required when connecting to very old NFS servers.

noexec

Prevents execution of binaries on mounted file systems. This is useful if the system is mounting a non-Linux file system containing incompatible binaries.

nosuid

Disables the **set-user-identifier** and **set-group-identifier** bits. This prevents remote users from gaining higher privileges by running a **setuid** program.

port=num

Specifies the numeric value of the NFS server port. If *num* is **0** (the default value), then **mount** queries the **rpcbind** service on the remote host for the port number to use. If the NFS service on the remote host is not registered with its **rpcbind** service, the standard NFS port number of TCP 2049 is used instead.

rsiZe=num and wsiZe=num

These options set the maximum number of bytes to be transferred in a single NFS read or write operation.

There is no fixed default value for **rsiZe** and **wsiZe**. By default, NFS uses the largest possible value that both the server and the client support. In Red Hat Enterprise Linux 9, the client and server maximum is 1,048,576 bytes. For more details, see the [What are the default and maximum values for rsiZe and wsiZe with NFS mounts?](#) KBase article.

sec=flavors

Security flavors to use for accessing files on the mounted export. The *flavors* value is a colon-separated list of one or more security flavors.

By default, the client attempts to find a security flavor that both the client and the server support. If the server does not support any of the selected flavors, the mount operation fails.

Available flavors:

- **sec=sys** uses local UNIX UIDs and GIDs. These use **AUTH_SYS** to authenticate NFS operations.
- **sec=krb5** uses Kerberos V5 instead of local UNIX UIDs and GIDs to authenticate users.
- **sec=krb5i** uses Kerberos V5 for user authentication and performs integrity checking of NFS operations using secure checksums to prevent data tampering.
- **sec=krb5p** uses Kerberos V5 for user authentication, integrity checking, and encrypts NFS traffic to prevent traffic sniffing. This is the most secure setting, but it also involves the most performance overhead.

tcp

Instructs the NFS mount to use the TCP protocol.

Additional resources

- **mount(8)** man page
- **nfs(5)** man page

3.9. STORING USER SETTINGS OVER NFS

If you use GNOME on a system with NFS home directories, you must set the **keyfile** back end for the **dconf** database. Otherwise, **dconf** might not work correctly. With this configuration, **dconf** stores settings in the `~/.config/dconf-keyfile/user` file.

Procedure

1. Ensure that the **glib2-fam** package is installed on the system:

```
# dnf install glib2-fam
```

Without this package, notifications on configuration changes made on remote machines are not displayed properly.

2. Create or edit the `/etc/dconf/profile/user` file on every client.
3. At the very beginning of the `/etc/dconf/profile/user` file, add the following line:

```
service-db:keyfile/user
```

4. Users must log out and log back in.
dconf polls the **keyfile** back end to determine whether updates have been made, so settings might not be updated immediately.

3.10. ADDITIONAL RESOURCES

- [The Linux NFS wiki](#)
- [Mounting NFS shares persistently](#).
- [Mounting NFS shares on demand](#).

CHAPTER 4. EXPORTING NFS SHARES

As a system administrator, you can use the NFS server to share a directory on your system over network.

4.1. INTRODUCTION TO NFS

This section explains the basic concepts of the NFS service.

A Network File System (NFS) allows remote hosts to mount file systems over a network and interact with those file systems as though they are mounted locally. This enables you to consolidate resources onto centralized servers on the network.

The NFS server refers to the **/etc/exports** configuration file to determine whether the client is allowed to access any exported file systems. Once verified, all file and directory operations are available to the user.

4.2. SUPPORTED NFS VERSIONS

This section lists versions of NFS supported in Red Hat Enterprise Linux and their features.

Currently, Red Hat Enterprise Linux 9 supports the following major versions of NFS:

- NFS version 3 (NFSv3) supports safe asynchronous writes and is more robust at error handling than the previous NFSv2; it also supports 64-bit file sizes and offsets, allowing clients to access more than 2 GB of file data.
- NFS version 4 (NFSv4) works through firewalls and on the Internet, no longer requires an **rpcbind** service, supports Access Control Lists (ACLs), and utilizes stateful operations.

NFS version 2 (NFSv2) is no longer supported by Red Hat.

Default NFS version

The default NFS version in Red Hat Enterprise Linux 9 is 4.2. NFS clients attempt to mount using NFSv4.2 by default, and fall back to NFSv4.1 when the server does not support NFSv4.2. The mount later falls back to NFSv4.0 and then to NFSv3.

Features of minor NFS versions

Following are the features of NFSv4.2 in Red Hat Enterprise Linux 9:

Server-side copy

Enables the NFS client to efficiently copy data without wasting network resources using the **copy_file_range()** system call.

Sparse files

Enables files to have one or more *holes*, which are unallocated or uninitialized data blocks consisting only of zeroes. The **lseek()** operation in NFSv4.2 supports **seek_hole()** and **seek_data()**, which enables applications to map out the location of holes in the sparse file.

Space reservation

Permits storage servers to reserve free space, which prohibits servers to run out of space. NFSv4.2 supports the **allocate()** operation to reserve space, the **deallocate()** operation to unreserve space, and the **fallocate()** operation to preallocate or deallocate space in a file.

Labeled NFS

Enforces data access rights and enables SELinux labels between a client and a server for individual files on an NFS file system.

Layout enhancements

Provides the **layoutstats()** operation, which enables some Parallel NFS (pNFS) servers to collect better performance statistics.

Following are the features of NFSv4.1:

- Enhances performance and security of network, and also includes client-side support for pNFS.
- No longer requires a separate TCP connection for callbacks, which allows an NFS server to grant delegations even when it cannot contact the client: for example, when NAT or a firewall interferes.
- Provides exactly once semantics (except for reboot operations), preventing a previous issue whereby certain operations sometimes returned an inaccurate result if a reply was lost and the operation was sent twice.

4.3. THE TCP AND UDP PROTOCOLS IN NFSV3 AND NFSV4

NFSv4 requires the Transmission Control Protocol (TCP) running over an IP network.

NFSv3 could also use the User Datagram Protocol (UDP) in earlier Red Hat Enterprise Linux versions. In Red Hat Enterprise Linux 9, NFS over UDP is no longer supported. By default, UDP is disabled in the NFS server.

4.4. SERVICES REQUIRED BY NFS

This section lists system services that are required for running an NFS server or mounting NFS shares. Red Hat Enterprise Linux starts these services automatically.

Red Hat Enterprise Linux uses a combination of kernel-level support and service processes to provide NFS file sharing. All NFS versions rely on Remote Procedure Calls (RPC) between clients and servers. To share or mount NFS file systems, the following services work together depending on which version of NFS is implemented:

nfsd

The NFS server kernel module that services requests for shared NFS file systems.

rpcbind

Accepts port reservations from local RPC services. These ports are then made available (or advertised) so the corresponding remote RPC services can access them. The **rpcbind** service responds to requests for RPC services and sets up connections to the requested RPC service. This is not used with NFSv4.

rpc.mountd

This process is used by an NFS server to process **MOUNT** requests from NFSv3 clients. It checks that the requested NFS share is currently exported by the NFS server, and that the client is allowed to access it. If the mount request is allowed, the **nfs-mountd** service replies with a Success status and provides the File-Handle for this NFS share back to the NFS client.

rpc.nfsd

This process enables explicit NFS versions and protocols the server advertises to be defined. It works with the Linux kernel to meet the dynamic demands of NFS clients, such as providing server threads each time an NFS client connects. This process corresponds to the **nfs-server** service.

lockd

This is a kernel thread that runs on both clients and servers. It implements the Network Lock Manager (NLM) protocol, which enables NFSv3 clients to lock files on the server. It is started automatically whenever the NFS server is run and whenever an NFS file system is mounted.

rpc.statd

This process implements the Network Status Monitor (NSM) RPC protocol, which notifies NFS clients when an NFS server is restarted without being gracefully brought down. The **rpc-statd** service is started automatically by the **nfs-server** service, and does not require user configuration. This is not used with NFSv4.

rpc.rquotad

This process provides user quota information for remote users. The **rpc-rquotad** service, which is provided by the **quota-rpc** package, has to be started by user when the **nfs-server** is started.

rpc.idmapd

This process provides NFSv4 client and server upcalls, which map between on-the-wire NFSv4 names (strings in the form of **user@domain**) and local UIDs and GIDs. For **idmapd** to function with NFSv4, the **/etc/idmapd.conf** file must be configured. At a minimum, the **Domain** parameter should be specified, which defines the NFSv4 mapping domain. If the NFSv4 mapping domain is the same as the DNS domain name, this parameter can be skipped. The client and server must agree on the NFSv4 mapping domain for ID mapping to function properly.

Only the NFSv4 server uses **rpc.idmapd**, which is started by the **nfs-idmapd** service. The NFSv4 client uses the keyring-based **nfsidmap** utility, which is called by the kernel on-demand to perform ID mapping. If there is a problem with **nfsidmap**, the client falls back to using **rpc.idmapd**.

The RPC services with NFSv4

The mounting and locking protocols have been incorporated into the NFSv4 protocol. The server also listens on the well-known TCP port 2049. As such, NFSv4 does not need to interact with **rpcbind**, **lockd**, and **rpc-statd** services. The **nfs-mountd** service is still required on the NFS server to set up the exports, but is not involved in any over-the-wire operations.

Additional resources

- [Configuring an NFSv4 only server without **rpcbind**](#).

4.5. NFS HOST NAME FORMATS

This section describes different formats that you can use to specify a host when mounting or exporting an NFS share.

You can specify the host in the following formats:

Single machine

Either of the following:

- A fully-qualified domain name (that can be resolved by the server)
- Host name (that can be resolved by the server)
- An IP address.

IP networks

Either of the following formats is valid:

- **a.b.c.d/z**, where **a.b.c.d** is the network and **z** is the number of bits in the netmask; for

example **192.168.0.0/24**.

- **a.b.c.d/netmask**, where **a.b.c.d** is the network and **netmask** is the netmask; for example, **192.168.100.8/255.255.255.0**.

Netgroups

The **@group-name** format, where **group-name** is the NIS netgroup name.

4.6. NFS SERVER CONFIGURATION

This section describes the syntax and options of two ways to configure exports on an NFS server:

- Manually editing the **/etc/exports** configuration file
- Using the **exportfs** utility on the command line

4.6.1. The **/etc/exports** configuration file

The **/etc/exports** file controls which file systems are exported to remote hosts and specifies options. It follows the following syntax rules:

- Blank lines are ignored.
- To add a comment, start a line with the hash mark (**#**).
- You can wrap long lines with a backslash (****).
- Each exported file system should be on its own individual line.
- Any lists of authorized hosts placed after an exported file system must be separated by space characters.
- Options for each of the hosts must be placed in parentheses directly after the host identifier, without any spaces separating the host and the first parenthesis.

Export entry

Each entry for an exported file system has the following structure:

```
export host(options)
```

It is also possible to specify multiple hosts, along with specific options for each host. To do so, list them on the same line as a space-delimited list, with each host name followed by its respective options (in parentheses), as in:

```
export host1(options1) host2(options2) host3(options3)
```

In this structure:

export

The directory being exported

host

The host or network to which the export is being shared

options

The options to be used for host

Example 4.1. A simple `/etc/exports` file

In its simplest form, the `/etc/exports` file only specifies the exported directory and the hosts permitted to access it:

```
/exported/directory bob.example.com
```

Here, **bob.example.com** can mount `/exported/directory/` from the NFS server. Because no options are specified in this example, NFS uses default options.

IMPORTANT

The format of the `/etc/exports` file is very precise, particularly in regards to use of the space character. Remember to always separate exported file systems from hosts and hosts from one another with a space character. However, there should be no other space characters in the file except on comment lines.

For example, the following two lines do not mean the same thing:

```
/home bob.example.com(rw)
/home bob.example.com (rw)
```

The first line allows only users from **bob.example.com** read and write access to the `/home` directory. The second line allows users from **bob.example.com** to mount the directory as read-only (the default), while the rest of the world can mount it read/write.

Default options

The default options for an export entry are:

ro

The exported file system is read-only. Remote hosts cannot change the data shared on the file system. To allow hosts to make changes to the file system (that is, read and write), specify the `rw` option.

sync

The NFS server will not reply to requests before changes made by previous requests are written to disk. To enable asynchronous writes instead, specify the option **async**.

wdelay

The NFS server will delay writing to the disk if it suspects another write request is imminent. This can improve performance as it reduces the number of times the disk must be accessed by separate write commands, thereby reducing write overhead. To disable this, specify the **no_wdelay** option, which is available only if the default `sync` option is also specified.

root_squash

This prevents root users connected remotely (as opposed to locally) from having root privileges; instead, the NFS server assigns them the user ID **nobody**. This effectively "squashes" the power of the remote root user to the lowest local user, preventing possible unauthorized writes on the remote server. To disable root squashing, specify the **no_root_squash** option.

To squash every remote user (including root), use the **all_squash** option. To specify the user and group IDs that the NFS server should assign to remote users from a particular host, use the **anonuid** and **anongid** options, respectively, as in:


```
export host(anonuid=uid,anongid=gid)
```

Here, *uid* and *gid* are user ID number and group ID number, respectively. The **anonuid** and **anongid** options enable you to create a special user and group account for remote NFS users to share.

By default, access control lists (ACLs) are supported by NFS under Red Hat Enterprise Linux. To disable this feature, specify the **no_acl** option when exporting the file system.

Default and overridden options

Each default for every exported file system must be explicitly overridden. For example, if the **rw** option is not specified, then the exported file system is shared as read-only. The following is a sample line from **/etc/exports** which overrides two default options:

```
/another/exported/directory 192.168.0.3(rw,async)
```

In this example, **192.168.0.3** can mount **/another/exported/directory/** read and write, and all writes to disk are asynchronous.

4.6.2. The **exportfs** utility

The **exportfs** utility enables the root user to selectively export or unexport directories without restarting the NFS service. When given the proper options, the **exportfs** utility writes the exported file systems to **/var/lib/nfs/xtab**. Because the **nfs-mountd** service refers to the **xtab** file when deciding access privileges to a file system, changes to the list of exported file systems take effect immediately.

Common **exportfs** options

The following is a list of commonly-used options available for **exportfs**:

-r

Causes all directories listed in **/etc/exports** to be exported by constructing a new export list in **/var/lib/nfs/etab**. This option effectively refreshes the export list with any changes made to **/etc/exports**.

-a

Causes all directories to be exported or unexported, depending on what other options are passed to **exportfs**. If no other options are specified, **exportfs** exports all file systems specified in **/etc/exports**.

-o file-systems

Specifies directories to be exported that are not listed in **/etc/exports**. Replace *file-systems* with additional file systems to be exported. These file systems must be formatted in the same way they are specified in **/etc/exports**. This option is often used to test an exported file system before adding it permanently to the list of exported file systems.

-i

Ignores **/etc/exports**; only options given from the command line are used to define exported file systems.

-u

Unexports all shared directories. The command **exportfs -ua** suspends NFS file sharing while keeping all NFS services up. To re-enable NFS sharing, use **exportfs -r**.

-v

Verbose operation, where the file systems being exported or unexported are displayed in greater detail when the **exportfs** command is executed.

If no options are passed to the **exportfs** utility, it displays a list of currently exported file systems.

Additional resources

- [NFS host name formats](#).

4.7. NFS AND RPCBIND

The **rpcbind** service maps Remote Procedure Call (RPC) services to the ports on which they listen. RPC processes notify **rpcbind** when they start, registering the ports they are listening on and the RPC program numbers they expect to serve. The client system then contacts **rpcbind** on the server with a particular RPC program number. The **rpcbind** service redirects the client to the proper port number so it can communicate with the requested service.

The Network File System Version 3 (NFSv3) requires the **rpcbind** service.

Because RPC-based services rely on **rpcbind** to make all connections with incoming client requests, **rpcbind** must be available before any of these services start.

Access control rules for **rpcbind** affect all RPC-based services. Alternatively, it is possible to specify access control rules for each of the NFS RPC daemons.

Additional resources

- **rpc.mountd(8)** man page
- **rpc.statd(8)** man page

4.8. INSTALLING NFS

This procedure installs all packages necessary to mount or export NFS shares.

Procedure

- Install the **nfs-utils** package:

```
# dnf install nfs-utils
```

4.9. STARTING THE NFS SERVER

This procedure describes how to start the NFS server, which is required to export NFS shares.

Prerequisites

- For servers that support NFSv3 connections, the **rpcbind** service must be running. To verify that **rpcbind** is active, use the following command:

```
$ systemctl status rpcbind
```

If the service is stopped, start and enable it:

```
$ systemctl enable --now rpcbind
```

Procedure

- To start the NFS server and enable it to start automatically at boot, use the following command:

```
# systemctl enable --now nfs-server
```

Additional resources

- [Configuring an NFSv4-only server.](#)

4.10. TROUBLESHOOTING NFS AND RPCBIND

Because the **rpcbind** service provides coordination between RPC services and the port numbers used to communicate with them, it is useful to view the status of current RPC services using **rpcbind** when troubleshooting. The **rpcinfo** utility shows each RPC-based service with port numbers, an RPC program number, a version number, and an IP protocol type (TCP or UDP).

Procedure

1. To make sure the proper NFS RPC-based services are enabled for **rpcbind**, use the following command:

```
# rpcinfo -p
```

Example 4.2. rpcinfo -p command output

The following is sample output from this command:

```
program vers proto  port  service
100000  4  tcp   111  portmapper
100000  3  tcp   111  portmapper
100000  2  tcp   111  portmapper
100000  4  udp   111  portmapper
100000  3  udp   111  portmapper
100000  2  udp   111  portmapper
100005  1  udp  20048 mountd
100005  1  tcp  20048 mountd
100005  2  udp  20048 mountd
100005  2  tcp  20048 mountd
100005  3  udp  20048 mountd
100005  3  tcp  20048 mountd
100024  1  udp  37769 status
100024  1  tcp  49349 status
100003  3  tcp   2049 nfs
100003  4  tcp   2049 nfs
100227  3  tcp   2049 nfs_acl
100021  1  udp  56691 nlockmgr
100021  3  udp  56691 nlockmgr
100021  4  udp  56691 nlockmgr
100021  1  tcp  46193 nlockmgr
100021  3  tcp  46193 nlockmgr
100021  4  tcp  46193 nlockmgr
```

If one of the NFS services does not start up correctly, **rpcbind** will be unable to map RPC requests from clients for that service to the correct port.

2. In many cases, if NFS is not present in **rpcinfo** output, restarting NFS causes the service to correctly register with **rpcbind** and begin working:

```
# systemctl restart nfs-server
```

Additional resources

- [Configuring an NFSv4-only server.](#)

4.11. CONFIGURING THE NFS SERVER TO RUN BEHIND A FIREWALL

NFS requires the **rpcbind** service, which dynamically assigns ports for RPC services and can cause issues for configuring firewall rules. The following sections describe how to configure NFS versions to work behind a firewall if you want to support:

- NFSv3
This includes any servers that support NFSv3:
 - NFSv3-only servers
 - Servers that support both NFSv3 and NFSv4
- NFSv4-only

4.11.1. Configuring the NFSv3-enabled server to run behind a firewall

The following procedure describes how to configure servers that support NFSv3 to run behind a firewall. This includes NFSv3-only servers and servers that support both NFSv3 and NFSv4.

Procedure

1. To allow clients to access NFS shares behind a firewall, configure the firewall by running the following commands on the NFS server:

```
firewall-cmd --permanent --add-service mountd
firewall-cmd --permanent --add-service rpc-bind
firewall-cmd --permanent --add-service nfs
```

2. Specify the ports to be used by the RPC service **nlockmgr** in the **/etc/nfs.conf** file as follows:

```
[lockd]

port=tcp-port-number
udp-port=udp-port-number
```

Alternatively, you can specify **nlm_tcpport** and **nlm_udpport** in the **/etc/modprobe.d/lockd.conf** file.

3. Open the specified ports in the firewall by running the following commands on the NFS server:

```
firewall-cmd --permanent --add-port=<lockd-tcp-port>/tcp
firewall-cmd --permanent --add-port=<lockd-udp-port>/udp
```

4. Add static ports for **rpc.statd** by editing the **[statd]** section of the **/etc/nfs.conf** file as follows:

```
[statd]

port=port-number
```

5. Open the added ports in the firewall by running the following commands on the NFS server:

```
firewall-cmd --permanent --add-port=<statd-tcp-port>/tcp
firewall-cmd --permanent --add-port=<statd-udp-port>/udp
```

6. Reload the firewall configuration:

```
firewall-cmd --reload
```

7. Restart the **rpc-statd** service first, and then restart the **nfs-server** service:

```
# systemctl restart rpc-statd.service
# systemctl restart nfs-server.service
```

Alternatively, if you specified the **lockd** ports in the **/etc/modprobe.d/lockd.conf** file:

- a. Update the current values of **/proc/sys/fs/nfs/nlm_tcpport** and **/proc/sys/fs/nfs/nlm_udpport**:

```
# sysctl -w fs.nfs.nlm_tcpport=<tcp-port>
# sysctl -w fs.nfs.nlm_udpport=<udp-port>
```

- b. Restart the **rpc-statd** and **nfs-server** services:

```
# systemctl restart rpc-statd.service
# systemctl restart nfs-server.service
```

4.11.2. Configuring the NFSv4-only server to run behind a firewall

The following procedure describes how to configure the NFSv4-only server to run behind a firewall.

Procedure

1. To allow clients to access NFS shares behind a firewall, configure the firewall by running the following command on the NFS server:

```
firewall-cmd --permanent --add-service nfs
```

2. Reload the firewall configuration:

```
firewall-cmd --reload
```

3. Restart the **nfs-server**:

```
# systemctl restart nfs-server
```

4.11.3. Configuring an NFSv3 client to run behind a firewall

The procedure to configure an NFSv3 client to run behind a firewall is similar to the procedure to configure an NFSv3 server to run behind a firewall.

If the machine you are configuring is both an NFS client and an NFS server, follow the procedure described in [Configuring the NFSv3-enabled server to run behind a firewall](#).

The following procedure describes how to configure a machine that is an NFS client only to run behind a firewall.

Procedure

1. To allow the NFS server to perform callbacks to the NFS client when the client is behind a firewall, add the **rpc-bind** service to the firewall by running the following command on the NFS client:

```
firewall-cmd --permanent --add-service rpc-bind
```

2. Specify the ports to be used by the RPC service **nlockmgr** in the **/etc/nfs.conf** file as follows:

```
[lockd]

port=port-number
udp-port=udp-port-number
```

Alternatively, you can specify **nlm_tcpport** and **nlm_udpport** in the **/etc/modprobe.d/lockd.conf** file.

3. Open the specified ports in the firewall by running the following commands on the NFS client:

```
firewall-cmd --permanent --add-port=<lockd-tcp-port>/tcp
firewall-cmd --permanent --add-port=<lockd-udp-port>/udp
```

4. Add static ports for **rpc.statd** by editing the **[statd]** section of the **/etc/nfs.conf** file as follows:

```
[statd]

port=port-number
```

5. Open the added ports in the firewall by running the following commands on the NFS client:

```
firewall-cmd --permanent --add-port=<statd-tcp-port>/tcp
firewall-cmd --permanent --add-port=<statd-udp-port>/udp
```

6. Reload the firewall configuration:

```
firewall-cmd --reload
```

7. Restart the **rpc-statd** service:

```
# systemctl restart rpc-statd.service
```

Alternatively, if you specified the **lockd** ports in the **/etc/modprobe.d/lockd.conf** file:

- a. Update the current values of **/proc/sys/fs/nfs/nlm_tcpport** and **/proc/sys/fs/nfs/nlm_udpport**:

```
# sysctl -w fs.nfs.nlm_tcpport=<tcp-port>
# sysctl -w fs.nfs.nlm_udpport=<udp-port>
```

- b. Restart the **rpc-statd** service:

```
# systemctl restart rpc-statd.service
```

4.11.4. Configuring an NFSv4 client to run behind a firewall

Perform this procedure only if the client is using NFSv4.0. In that case, it is necessary to open a port for NFSv4.0 callbacks.

This procedure is not needed for NFSv4.1 or higher because in the later protocol versions the server performs callbacks on the same connection that was initiated by the client.

Procedure

1. To allow NFSv4.0 callbacks to pass through firewalls, set **/proc/sys/fs/nfs/nfs_callback_tcpport** and allow the server to connect to that port on the client as follows:

```
# echo "fs.nfs.nfs_callback_tcpport = <callback-port>" >/etc/sysctl.d/90-nfs-callback-
port.conf
# sysctl -p /etc/sysctl.d/90-nfs-callback-port.conf
```

2. Open the specified port in the firewall by running the following command on the NFS client:

```
firewall-cmd --permanent --add-port=<callback-port>/tcp
```

3. Reload the firewall configuration:

```
firewall-cmd --reload
```

4.12. EXPORTING RPC QUOTA THROUGH A FIREWALL

If you export a file system that uses disk quotas, you can use the quota Remote Procedure Call (RPC) service to provide disk quota data to NFS clients.

Procedure

1. Enable and start the **rpc-rquotad** service:

```
# systemctl enable --now rpc-rquotad
```

**NOTE**

The **rpc-rquotad** service is, if enabled, started automatically after starting the **nfs-server** service.

2. To make the quota RPC service accessible behind a firewall, the TCP (or UDP, if UDP is enabled) port 875 need to be open. The default port number is defined in the **/etc/services** file. You can override the default port number by appending **-p port-number** to the **RPCRQUOTADOPTS** variable in the **/etc/sysconfig/rpc-rquotad** file.
3. By default, remote hosts can only read quotas. If you want to allow clients to set quotas, append the **-S** option to the **RPCRQUOTADOPTS** variable in the **/etc/sysconfig/rpc-rquotad** file.
4. Restart **rpc-rquotad** for the changes in the **/etc/sysconfig/rpc-rquotad** file to take effect:

```
# systemctl restart rpc-rquotad
```

4.13. ENABLING NFS OVER RDMA (NFSORDMA)

In Red Hat Enterprise Linux 9, Remote direct memory access (RDMA) service on RDMA-capable hardware provides Network File System (NFS) protocol support for high-speed file transfer over the network.

Procedure

1. Install the **rdma-core** package:

```
# dnf install rdma-core
```

2. Verify the lines with **xprtrdma** and **svcrdma** are not commented out in the **/etc/rdma/modules/rdma.conf** file:

```
# NFS over RDMA client support
xprtrdma
# NFS over RDMA server support
svcrdma
```

3. On the NFS server, create directory **/mnt/nfsordma** and export it to **/etc/exports**:

```
# mkdir /mnt/nfsordma
# echo "/mnt/nfsordma *(fsid=0,rw,async,insecure,no_root_squash)" >> /etc/exports
```

4. On the NFS client, mount the nfs-share with server IP address, for example, **172.31.0.186**:

```
# mount -o rdma,port=20049 172.31.0.186:/mnt/nfs-share /mnt/nfs
```

5. Restart the **nfs-server** service:

```
# systemctl restart nfs-server
```

Additional resources

- [The RFC 5667 standard](#)

4.14. ADDITIONAL RESOURCES

- [The Linux NFS wiki](#)

CHAPTER 5. CONFIGURING AN NFSV4-ONLY SERVER

As an NFS server administrator, you can configure the NFS server to support only NFSv4, which minimizes the number of open ports and running services on the system.

5.1. BENEFITS AND DRAWBACKS OF AN NFSV4-ONLY SERVER

This section explains the benefits and drawbacks of configuring the NFS server to only support NFSv4.

By default, the NFS server supports NFSv3 and NFSv4 connections in Red Hat Enterprise Linux 9. However, you can also configure NFS to support only NFS version 4.0 and later. This minimizes the number of open ports and running services on the system, because NFSv4 does not require the **rpcbind** service to listen on the network.

When your NFS server is configured as NFSv4-only, clients attempting to mount shares using NFSv3 fail with an error like the following:

```
Requested NFS version or transport protocol is not supported.
```

Optionally, you can also disable listening for the **RPCBIND**, **MOUNT**, and **NSM** protocol calls, which are not necessary in the NFSv4-only case.

The effects of disabling these additional options are:

- Clients that attempt to mount shares from your server using NFSv3 become unresponsive.
- The NFS server itself is unable to mount NFSv3 file systems.

5.2. CONFIGURING THE NFS SERVER TO SUPPORT ONLY NFSV4

This procedure describes how to configure your NFS server to support only NFS version 4.0 and later.

Procedure

1. Disable NFSv3 by adding the following lines to the **[nfsd]** section of the **/etc/nfs.conf** configuration file:

```
[nfsd]
vers3=no
```

2. Optionally, disable listening for the **RPCBIND**, **MOUNT**, and **NSM** protocol calls, which are not necessary in the NFSv4-only case. Disable related services:

```
# systemctl mask --now rpc-statd.service rpcbind.service rpcbind.socket
```

3. Restart the NFS server:

```
# systemctl restart nfs-server
```

The changes take effect as soon as you start or restart the NFS server.

5.3. VERIFYING THE NFSV4-ONLY CONFIGURATION

You can verify that your NFS server is configured in the NFSv4-only mode by using the **netstat** utility.

Procedure

- Use the **netstat** utility to list services listening on the TCP and UDP protocols:

```
# netstat --listening --tcp --udp
```

Example 5.1. Output on an NFSv4-only server

The following is an example **netstat** output on an NFSv4-only server; listening for **RPCBIND** and **NSM** is also disabled.

Note that the **mountd** is still listening, but only because internal communication between the **nfsd** process and **mountd** is required in order to do local authentication, therefore the port can remain closed.

```
# netstat --listening --tcp --udp
```

Active Internet connections (only servers)

Proto	Recv-Q	Send-Q	Local Address	Foreign Address	State
tcp	0	0	0.0.0.0:nfs	0.0.0.0:*	LISTEN
tcp	0	0	0.0.0.0:mountd	0.0.0.0:*	LISTEN
tcp	0	0	0.0.0.0:ssh	0.0.0.0:*	LISTEN
tcp6	0	0	:::nfs	:::*	LISTEN
tcp6	0	0	:::mountd	:::*	LISTEN
tcp6	0	0	:::ssh	:::*	LISTEN
udp	0	0	0.0.0.0:mountd	0.0.0.0:*	
udp	0	0	0.0.0.0:mountd	0.0.0.0:*	
udp	0	0	localhost:323	0.0.0.0:*	
udp6	0	0	:::mountd	:::*	
udp6	0	0	:::mountd	:::*	
udp6	0	0	localhost:323	:::*	

Example 5.2. Output before configuring an NFSv4-only server

In comparison, the **netstat** output before configuring an NFSv4-only server includes the **sunrpc** service:

```
# netstat --listening --tcp --udp
```

Active Internet connections (only servers)

Proto	Recv-Q	Send-Q	Local Address	Foreign Address	State
tcp	0	0	0.0.0.0:ssh	0.0.0.0:*	LISTEN
tcp	0	0	0.0.0.0:40189	0.0.0.0:*	LISTEN
tcp	0	0	0.0.0.0:46813	0.0.0.0:*	LISTEN
tcp	0	0	0.0.0.0:nfs	0.0.0.0:*	LISTEN
tcp	0	0	0.0.0.0:sunrpc	0.0.0.0:*	LISTEN
tcp	0	0	0.0.0.0:mountd	0.0.0.0:*	LISTEN
tcp6	0	0	:::ssh	:::*	LISTEN
tcp6	0	0	:::51227	:::*	LISTEN

```
tcp6 0 0 [::]:nfs [::]:* LISTEN
tcp6 0 0 [::]:sunrpc [::]:* LISTEN
tcp6 0 0 [::]:mountd [::]:* LISTEN
tcp6 0 0 [::]:45043 [::]:* LISTEN
udp 0 0 localhost:1018 0.0.0.0:*
udp 0 0 localhost.locald:bootpc 0.0.0.0:*
udp 0 0 0.0.0.0:mountd 0.0.0.0:*
udp 0 0 0.0.0.0:46672 0.0.0.0:*
udp 0 0 0.0.0.0:sunrpc 0.0.0.0:*
udp 0 0 0.0.0.0:33494 0.0.0.0:*
udp6 0 0 [::]:33734 [::]:*
udp6 0 0 [::]:mountd [::]:*
udp6 0 0 [::]:sunrpc [::]:*
udp6 0 0 [::]:40243 [::]:*
```

CHAPTER 6. SECURING NFS

To minimize NFS security risks and protect data on the server, consider the following sections when exporting NFS file systems on a server or mounting them on a client.

6.1. NFS SECURITY WITH AUTH_SYS AND EXPORT CONTROLS

NFS provides the following traditional options in order to control access to exported files:

- The server restricts which hosts are allowed to mount which file systems either by IP address or by host name.
- The server enforces file system permissions for users on NFS clients in the same way it does for local users. Traditionally, NFS does this using the **AUTH_SYS** call message (also called **AUTH_UNIX**), which relies on the client to state the UID and GIDs of the user. Be aware that this means that a malicious or misconfigured client might easily get this wrong and allow a user access to files that it should not.

To limit the potential risks, administrators often limits the access to read-only or squash user permissions to a common user and group ID. Unfortunately, these solutions prevent the NFS share from being used in the way it was originally intended.

Additionally, if an attacker gains control of the DNS server used by the system exporting the NFS file system, they can point the system associated with a particular hostname or fully qualified domain name to an unauthorized machine. At this point, the unauthorized machine *is* the system permitted to mount the NFS share, because no username or password information is exchanged to provide additional security for the NFS mount.

Wildcards should be used sparingly when exporting directories through NFS, as it is possible for the scope of the wildcard to encompass more systems than intended.

Additional resources

- To secure NFS and **rpcbind**, use, for example, **nftables** and **firewalld**.
- **nft(8)** man page
- **firewalld-cmd(1)** man page

6.2. NFS SECURITY WITH AUTH_GSS

All version of NFS support RPCSEC_GSS and the Kerberos mechanism.

Unlike AUTH_SYS, with the RPCSEC_GSS Kerberos mechanism, the server does not depend on the client to correctly represent which user is accessing the file. Instead, cryptography is used to authenticate users to the server, which prevents a malicious client from impersonating a user without having that user's Kerberos credentials. Using the RPCSEC_GSS Kerberos mechanism is the most straightforward way to secure mounts because after configuring Kerberos, no additional setup is needed.

6.3. CONFIGURING AN NFS SERVER AND CLIENT TO USE KERBEROS

Kerberos is a network authentication system that allows clients and servers to authenticate to each other by using symmetric encryption and a trusted third party, the KDC. Red Hat recommends using Identity Management (IdM) for setting up Kerberos.

Prerequisites

- The Kerberos Key Distribution Centre (**KDC**) is installed and configured.

Procedure

1.
 - Create the **nfs/hostname.domain@REALM** principal on the NFS server side.
 - Create the **host/hostname.domain@REALM** principal on both the server and the client side.
 - Add the corresponding keys to keytabs for the client and server.
2. On the server side, use the **sec=** option to enable the wanted security flavors. To enable all security flavors as well as non-cryptographic mounts:

```
/export *(sec=sys:krb5:krb5i:krb5p)
```

Valid security flavors to use with the **sec=** option are:

- **sys**: no cryptographic protection, the default
 - **krb5**: authentication only
 - **krb5i**: integrity protection
 - uses Kerberos V5 for user authentication and performs integrity checking of NFS operations using secure checksums to prevent data tampering.
 - **krb5p**: privacy protection
 - uses Kerberos V5 for user authentication, integrity checking, and encrypts NFS traffic to prevent traffic sniffing. This is the most secure setting, but it also involves the most performance overhead.
3. On the client side, add **sec=krb5** (or **sec=krb5i**, or **sec=krb5p**, depending on the setup) to the mount options:

```
# mount -o sec=krb5 server:/export /mnt
```

Additional resources

- [Creating files as root on krb5-secured NFS](#) . Not recommended.
- **exports(5)** man page
- **nfs(5)** man page

6.4. NFSV4 SECURITY OPTIONS

NFSv4 includes ACL support based on the Microsoft Windows NT model, not the POSIX model, because of the Microsoft Windows NT model's features and wide deployment.

Another important security feature of NFSv4 is the removal of the use of the **MOUNT** protocol for mounting file systems. The **MOUNT** protocol presented a security risk because of the way the protocol processed file handles.

6.5. FILE PERMISSIONS ON MOUNTED NFS EXPORTS

Once the NFS file system is mounted as either read or read and write by a remote host, the only protection each shared file has is its permissions. If two users that share the same user ID value mount the same NFS file system on different client systems, they can modify each others' files. Additionally, anyone logged in as root on the client system can use the **su -** command to access any files with the NFS share.

By default, access control lists (ACLs) are supported by NFS under Red Hat Enterprise Linux. Red Hat recommends to keep this feature enabled.

By default, NFS uses *root squashing* when exporting a file system. This sets the user ID of anyone accessing the NFS share as the root user on their local machine to **nobody**. Root squashing is controlled by the default option **root_squash**; for more information about this option, see [NFS server configuration](#).

When exporting an NFS share as read-only, consider using the **all_squash** option. This option makes every user accessing the exported file system take the user ID of the **nobody** user.

CHAPTER 7. ENABLING PNFS SCSI LAYOUTS IN NFS

You can configure the NFS server and client to use the pNFS SCSI layout for accessing data. pNFS SCSI is beneficial in use cases that involve longer-duration single-client access to a file.

Prerequisites

- Both the client and the server must be able to send SCSI commands to the same block device. That is, the block device must be on a shared SCSI bus.
- The block device must contain an XFS file system.
- The SCSI device must support SCSI Persistent Reservations as described in the SCSI-3 Primary Commands specification.

7.1. THE PNFS TECHNOLOGY

The pNFS architecture improves the scalability of NFS. When a server implements pNFS, the client is able to access data through multiple servers concurrently. This can lead to performance improvements.

pNFS supports the following storage protocols or layouts on RHEL:

- Files
- Flexfiles
- SCSI

7.2. PNFS SCSI LAYOUTS

The SCSI layout builds on the work of pNFS block layouts. The layout is defined across SCSI devices. It contains a sequential series of fixed-size blocks as logical units (LUs) that must be capable of supporting SCSI persistent reservations. The LU devices are identified by their SCSI device identification.

pNFS SCSI performs well in use cases that involve longer-duration single-client access to a file. An example might be a mail server or a virtual machine housing a cluster.

Operations between the client and the server

When an NFS client reads from a file or writes to it, the client performs a **LAYOUTGET** operation. The server responds with the location of the file on the SCSI device. The client might need to perform an additional operation of **GETDEVICEINFO** to determine which SCSI device to use. If these operations work correctly, the client can issue I/O requests directly to the SCSI device instead of sending **READ** and **WRITE** operations to the server.

Errors or contention between clients might cause the server to recall layouts or not issue them to the clients. In those cases, the clients fall back to issuing **READ** and **WRITE** operations to the server instead of sending I/O requests directly to the SCSI device.

To monitor the operations, see [Monitoring pNFS SCSI layouts functionality](#).

Device reservations

pNFS SCSI handles fencing through the assignment of reservations. Before the server issues layouts to clients, it reserves the SCSI device to ensure that only registered clients may access the device. If a client can issue commands to that SCSI device but is not registered with the device, many operations

from the client on that device fail. For example, the **blkid** command on the client fails to show the UUID of the XFS file system if the server has not given a layout for that device to the client.

The server does not remove its own persistent reservation. This protects the data within the file system on the device across restarts of clients and servers. In order to repurpose the SCSI device, you might need to manually remove the persistent reservation on the NFS server.

7.3. CHECKING FOR A SCSI DEVICE COMPATIBLE WITH PNFS

This procedure checks if a SCSI device supports the pNFS SCSI layout.

Prerequisites

- Install the **sg3_utils** package:

```
# dnf install sg3_utils
```

Procedure

- On both the server and client, check for the proper SCSI device support:

```
# sg_persist --in --report-capabilities --verbose path-to-scsi-device
```

Ensure that the *Persist Through Power Loss Active* (**PTPL_A**) bit is set.

Example 7.1. A SCSI device that supports pNFS SCSI

The following is an example of **sg_persist** output for a SCSI device that supports pNFS SCSI. The **PTPL_A** bit reports **1**.

```
inquiry cdb: 12 00 00 00 24 00
Persistent Reservation In cmd: 5e 02 00 00 00 00 00 20 00 00
LIO-ORG block11 4.0
Peripheral device type: disk
Report capabilities response:
Compatible Reservation Handling(CRH): 1
Specify Initiator Ports Capable(SIP_C): 1
All Target Ports Capable(ATP_C): 1
Persist Through Power Loss Capable(PTPL_C): 1
Type Mask Valid(TMV): 1
Allow Commands: 1
Persist Through Power Loss Active(PTPL_A): 1
Support indicated in Type mask:
Write Exclusive, all registrants: 1
Exclusive Access, registrants only: 1
Write Exclusive, registrants only: 1
Exclusive Access: 1
Write Exclusive: 1
Exclusive Access, all registrants: 1
```

Additional resources

- **sg_persist(8)** man page

7.4. SETTING UP PNFS SCSI ON THE SERVER

This procedure configures an NFS server to export a pNFS SCSI layout.

Procedure

1. On the server, mount the XFS file system created on the SCSI device.
2. Configure the NFS server to export NFS version 4.1 or higher. Set the following option in the **[nfsd]** section of the **/etc/nfs.conf** file:

```
[nfsd]
vers4.1=y
```

3. Configure the NFS server to export the XFS file system over NFS with the **pnfs** option:

Example 7.2. An entry in /etc/exports to export pNFS SCSI

The following entry in the **/etc/exports** configuration file exports the file system mounted at **/exported/directory/** to the **allowed.example.com** client as a pNFS SCSI layout:

```
/exported/directory allowed.example.com(pnfs)
```



NOTE

The exported file system must be created on the whole block device, not only on a partition.

Additional resources

- [Exporting NFS shares](#).

7.5. SETTING UP PNFS SCSI ON THE CLIENT

This procedure configures an NFS client to mount a pNFS SCSI layout.

Prerequisites

- The NFS server is configured to export an XFS file system over pNFS SCSI. See [Setting up pNFS SCSI on the server](#).

Procedure

- On the client, mount the exported XFS file system using NFS version 4.1 or higher:

```
# mount -t nfs -o nfsvers=4.1 host:/remote/export /local/directory
```

Do not mount the XFS file system directly without NFS.

Additional resources

- [Mounting NFS shares.](#)

7.6. RELEASING THE PNFS SCSI RESERVATION ON THE SERVER

This procedure releases the persistent reservation that an NFS server holds on a SCSI device. This enables you to repurpose the SCSI device when you no longer need to export pNFS SCSI.

You must remove the reservation from the server. It cannot be removed from a different IT Nexus.

Prerequisites

- Install the **sg3_utils** package:

```
# dnf install sg3_utils
```

Procedure

1. Query an existing reservation on the server:

```
# sg_persist --read-reservation path-to-scsi-device
```

Example 7.3. Querying a reservation on /dev/sda

```
# sg_persist --read-reservation /dev/sda

LIO-ORG  block_1      4.0
Peripheral device type: disk
PR generation=0x8, Reservation follows:
Key=0x1000000000000000
scope: LU_SCOPE, type: Exclusive Access, registrants only
```

2. Remove the existing registration on the server:

```
# sg_persist --out \
    --release \
    --param-rk=reservation-key \
    --prout-type=6 \
    path-to-scsi-device
```

Example 7.4. Removing a reservation on /dev/sda

```
# sg_persist --out \
    --release \
    --param-rk=0x1000000000000000 \
    --prout-type=6 \
    /dev/sda

LIO-ORG  block_1      4.0
Peripheral device type: disk
```

■

Additional resources

- **sg_persist(8)** man page

CHAPTER 8. MONITORING PNFS SCSI LAYOUTS FUNCTIONALITY

You can monitor if the pNFS client and server exchange proper pNFS SCSI operations or if they fall back on regular NFS operations.

Prerequisites

- A pNFS SCSI client and server are configured.

8.1. CHECKING PNFS SCSI OPERATIONS FROM THE SERVER USING NFSSTAT

This procedure uses the **nfsstat** utility to monitor pNFS SCSI operations from the server.

Procedure

1. Monitor the operations serviced from the server:

```
# watch --differences \
    "nfsstat --server | egrep --after-context=1 read\|write\|layout"

Every 2.0s: nfsstat --server | egrep --after-context=1 read\|write\|layout

putrootfh  read      readdir  readlink  remove  rename
2         0% 0      0% 1      0% 0      0% 0      0% 0      0%
--
setctidconf verify  write      rellockowner bc_ctl  bind_conn
0         0% 0      0% 0      0% 0      0% 0      0% 0
--
getdevlist layoutcommit layoutget  layoutreturn secinfoonam sequence
0         0% 29     1% 49     1% 5      0% 0      0% 2435  86%
```

2. The client and server use pNFS SCSI operations when:

- The **layoutget**, **layoutreturn**, and **layoutcommit** counters increment. This means that the server is serving layouts.
- The server **read** and **write** counters do not increment. This means that the clients are performing I/O requests directly to the SCSI devices.

8.2. CHECKING PNFS SCSI OPERATIONS FROM THE CLIENT USING MOUNTSTATS

This procedure uses the **/proc/self/mountstats** file to monitor pNFS SCSI operations from the client.

Procedure

1. List the per-mount operation counters:

```
# cat /proc/self/mountstats \
    | awk /scsi_lun_0/,/^$/ \
```

```
| egrep device\|READ\|WRITE\|LAYOUT
```

```
device 192.168.122.73:/exports/scsi_lun_0 mounted on /mnt/rhel7/scsi_lun_0 with fstype  
nfs4 statvers=1.1
```

```
  nfsv4:
```

```
bm0=0xfdffbfff,bm1=0x40f9be3e,bm2=0x803,acl=0x3,sessions,pnfs=LAYOUT_SCSI
```

```
    READ: 0 0 0 0 0 0 0 0
```

```
    WRITE: 0 0 0 0 0 0 0 0
```

```
    READLINK: 0 0 0 0 0 0 0 0
```

```
    READDIR: 0 0 0 0 0 0 0 0
```

```
    LAYOUTGET: 49 49 0 11172 9604 2 19448 19454
```

```
    LAYOUTCOMMIT: 28 28 0 7776 4808 0 24719 24722
```

```
    LAYOUTRETURN: 0 0 0 0 0 0 0 0
```

```
    LAYOUTSTATS: 0 0 0 0 0 0 0 0
```

2. In the results:

- The **LAYOUT** statistics indicate requests where the client and server use pNFS SCSI operations.
- The **READ** and **WRITE** statistics indicate requests where the client and server fall back to NFS operations.

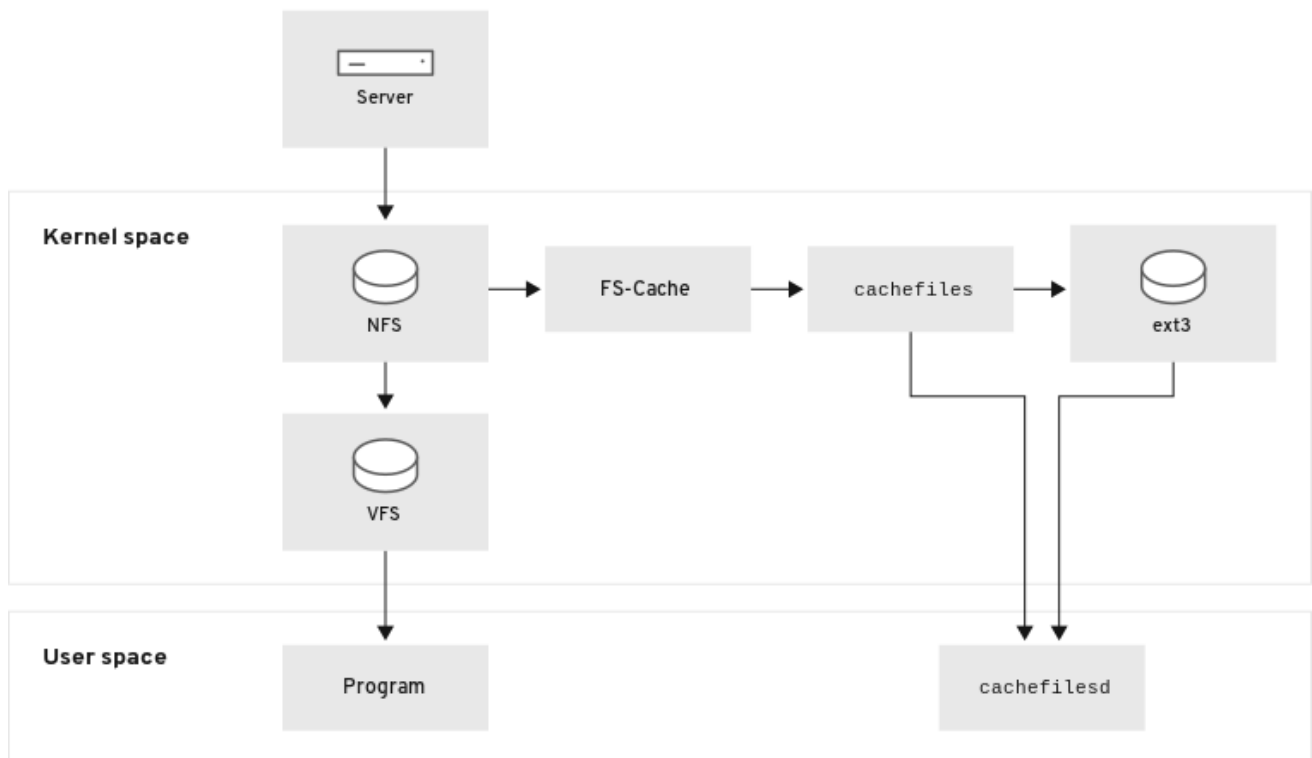
CHAPTER 9. GETTING STARTED WITH FS-CACHE

FS-Cache is a persistent local cache that file systems can use to take data retrieved from over the network and cache it on local disk. This helps minimize network traffic for users accessing data from a file system mounted over the network (for example, NFS).

9.1. OVERVIEW OF THE FS-CACHE

The following diagram is a high-level illustration of how FS-Cache works:

Figure 9.1. FS-Cache Overview



96_RHEL_0720

FS-Cache is designed to be as transparent as possible to the users and administrators of a system. Unlike **cachefs** on Solaris, FS-Cache allows a file system on a server to interact directly with a client's local cache without creating an overmounted file system. With NFS, a mount option instructs the client to mount the NFS share with FS-cache enabled. The mount point will cause automatic upload for two kernel modules: **fscache** and **cachefiles**. The **cachefilesd** daemon communicates with the kernel modules to implement the cache.

FS-Cache does not alter the basic operation of a file system that works over the network - it merely provides that file system with a persistent place in which it can cache data. For instance, a client can still mount an NFS share whether or not FS-Cache is enabled. In addition, cached NFS can handle files that will not fit into the cache (whether individually or collectively) as files can be partially cached and do not have to be read completely up front. FS-Cache also hides all I/O errors that occur in the cache from the client file system driver.

To provide caching services, FS-Cache needs a *cache back end*. A cache back end is a storage driver configured to provide caching services, which is **cachefiles**. In this case, FS-Cache requires a mounted block-based file system that supports **bmap** and extended attributes (e.g. ext3) as its cache back end.

File systems that support functionalities required by FS-Cache cache back end include the Red Hat Enterprise Linux 9 implementations of the following file systems:

- ext3 (with extended attributes enabled)
- ext4
- XFS

FS-Cache cannot arbitrarily cache any file system, whether through the network or otherwise: the shared file system's driver must be altered to allow interaction with FS-Cache, data storage/retrieval, and metadata setup and validation. FS-Cache needs *indexing keys* and *coherency data* from the cached file system to support persistence: indexing keys to match file system objects to cache objects, and coherency data to determine whether the cache objects are still valid.



NOTE

In Red Hat Enterprise Linux 9, the **cachefilesd** package is not installed by default and needs to be installed manually.

9.2. PERFORMANCE GUARANTEE

FS-Cache does *not* guarantee increased performance. Using a cache incurs a performance penalty: for example, cached NFS shares add disk accesses to cross-network lookups. While FS-Cache tries to be as asynchronous as possible, there are synchronous paths (e.g. reads) where this isn't possible.

For example, using FS-Cache to cache an NFS share between two computers over an otherwise unladen GigE network likely will not demonstrate any performance improvements on file access. Rather, NFS requests would be satisfied faster from server memory rather than from local disk.

The use of FS-Cache, therefore, is a *compromise* between various factors. If FS-Cache is being used to cache NFS traffic, for instance, it may slow the client down a little, but massively reduce the network and server loading by satisfying read requests locally without consuming network bandwidth.

9.3. SETTING UP A CACHE

Currently, Red Hat Enterprise Linux 9 only provides the **cachefiles** caching back end. The **cachefilesd** daemon initiates and manages **cachefiles**. The **/etc/cachefilesd.conf** file controls how **cachefiles** provides caching services.

The cache back end works by maintaining a certain amount of free space on the partition hosting the cache. It grows and shrinks the cache in response to other elements of the system using up free space, making it safe to use on the root file system (for example, on a laptop). FS-Cache sets defaults on this behavior, which can be configured via *cache cull limits*. For more information about configuring cache cull limits, see [Cache cull limits configuration](#).

This procedure shows how to set up a cache.

Prerequisites

- The **cachefilesd** package is installed and service has started successfully. To be sure the service is running, use the following command:

```
# systemctl start cachefilesd
# systemctl status cachefilesd
```


The status must be *active (running)*.

Procedure

1. Configure in a cache back end which directory to use as a cache, use the following parameter:

```
$ dir /path/to/cache
```

2. Typically, the cache back end directory is set in **/etc/cachefilesd.conf** as **/var/cache/fscache**, as in:

```
$ dir /var/cache/fscache
```

3. If you want to change the cache back end directory, the selinux context must be same as **/var/cache/fscache**:

```
# semanage fcontext -a -e /var/cache/fscache /path/to/cache
# restorecon -Rv /path/to/cache
```

4. Replace */path/to/cache* with the directory name while setting up cache.
5. If the given commands for setting selinux context did not work, use the following commands:

```
# semanage permissive -a cachefilesd_t
# semanage permissive -a cachefiles_kernel_t
```

FS-Cache will store the cache in the file system that hosts **/path/to/cache**. On a laptop, it is advisable to use the root file system (/) as the host file system, but for a desktop machine it would be more prudent to mount a disk partition specifically for the cache.

6. The host file system must support user-defined extended attributes; FS-Cache uses these attributes to store coherency maintenance information. To enable user-defined extended attributes for ext3 file systems (i.e. **device**), use:

```
# tune2fs -o user_xattr /dev/device
```

7. To enable extended attributes for a file system at the mount time, as an alternative, use the following command:

```
# mount /dev/device /path/to/cache -o user_xattr
```

8. Once the configuration file is in place, start up the **cachefilesd** service:

```
# systemctl start cachefilesd
```

9. To configure **cachefilesd** to start at boot time, execute the following command as root:

```
# systemctl enable cachefilesd
```

9.4. CACHE CULL LIMITS CONFIGURATION

The **cachefilesd** daemon works by caching remote data from shared file systems to free space on the

disk. This could potentially consume all available free space, which could be bad if the disk also housed the root partition. To control this, **cachefilesd** tries to maintain a certain amount of free space by discarding old objects (i.e. accessed less recently) from the cache. This behavior is known as *cache culling*.

Cache culling is done on the basis of the percentage of blocks and the percentage of files available in the underlying file system. There are settings in **/etc/cachefilesd.conf** which control six limits:

brun *N*% (percentage of blocks), frun *N*% (percentage of files)

If the amount of free space and the number of available files in the cache rises above both these limits, then culling is turned off.

bcull *N*% (percentage of blocks), fcull *N*% (percentage of files)

If the amount of available space or the number of files in the cache falls below either of these limits, then culling is started.

bstop *N*% (percentage of blocks), fstop *N*% (percentage of files)

If the amount of available space or the number of available files in the cache falls below either of these limits, then no further allocation of disk space or files is permitted until culling has raised things above these limits again.

The default value of **N** for each setting is as follows:

- **brun/frun** - 10%
- **bcull/fcull** - 7%
- **bstop/fstop** - 3%

When configuring these settings, the following must hold true:

- $0 \leq \mathbf{bstop} < \mathbf{bcull} < \mathbf{brun} < 100$
- $0 \leq \mathbf{fstop} < \mathbf{fcull} < \mathbf{frun} < 100$

These are the percentages of available space and available files and do not appear as 100 minus the percentage displayed by the **df** program.



IMPORTANT

Culling depends on both **bxxx** and **fxxx** pairs simultaneously; the user can not treat them separately.

9.5. RETRIEVING STATISTICAL INFORMATION FROM THE FSCACHE KERNEL MODULE

FS-Cache also keeps track of general statistical information. This procedure shows how to get this information.

Procedure

1. To view the statistical information about FS-Cache, use the following command:

```
# cat /proc/fs/fscache/stats
```

FS-Cache statistics includes information about decision points and object counters. For more information, see the following kernel document:

`/usr/share/doc/kernel-doc-4.18.0/Documentation/filesystems/caching/fscache.txt`

9.6. FS-CACHE REFERENCES

This section provides reference information for FS-Cache.

1. For more information about **cachefilesd** and how to configure it, see **man cachefilesd** and **man cachefilesd.conf**. The following kernel documents also provide additional information:
 - **`/usr/share/doc/cachefilesd/README`**
 - **`/usr/share/man/man5/cachefilesd.conf.5.gz`**
 - **`/usr/share/man/man8/cachefilesd.8.gz`**
2. For general information about FS-Cache, including details on its design constraints, available statistics, and capabilities, see the following kernel document:
`/usr/share/doc/kernel-doc-4.18.0/Documentation/filesystems/caching/fscache.txt`

CHAPTER 10. USING THE CACHE WITH NFS

NFS will not use the cache unless explicitly instructed. This paragraph shows how to configure an NFS mount by using FS-Cache.

Prerequisites

- The **cachefilesd** package is installed and running. To ensure it is running, use the following command:

```
# systemctl start cachefilesd
# systemctl status cachefilesd
```

The status must be *active (running)*.

- Mount NFS shares with the following option:

```
# mount nfs-share:/ /mount/point -o fsc
```

All access to files under **/mount/point** will go through the cache, unless the file is opened for direct I/O or writing. For more information, see [Cache limitations with NFS](#).

NFS indexes cache contents using NFS file handle, *not* the file name, which means hard-linked files share the cache correctly.

NFS versions 3, 4.0, 4.1 and 4.2 support caching. However, each version uses different branches for caching.

10.1. CONFIGURING NFS CACHE SHARING

There are several potential issues to do with NFS cache sharing. Because the cache is persistent, blocks of data in the cache are indexed on a sequence of four keys:

- Level 1: Server details
- Level 2: Some mount options; security type; FSID; uniquifier
- Level 3: File Handle
- Level 4: Page number in file

To avoid coherency management problems between superblocks, all NFS superblocks that require to cache the data have unique Level 2 keys. Normally, two NFS mounts with same source volume and options share a superblock, and therefore share the caching, even if they mount different directories within that volume.

This is an example how to configure cache sharing with different options.

Procedure

1. Mount NFS shares with the following commands:

```
mount home0:/disk0/fred /home/fred -o fsc
mount home0:/disk0/jim /home/jim -o fsc
```

Here, **/home/fred** and **/home/jim** likely share the superblock as they have the same options, especially if they come from the same volume/partition on the NFS server (**home0**).

- To not share the superblock, use the **mount** command with the following options:

```
mount home0:/disk0/fred /home/fred -o fsc,rsize=8192
mount home0:/disk0/jim /home/jim -o fsc,rsize=65536
```

In this case, **/home/fred** and **/home/jim** will not share the superblock as they have different network access parameters, which are part of the Level 2 key.

- To cache the contents of the two subtrees (**/home/fred1** and **/home/fred2**) *twice* with not sharing the superblock, use the following command:

```
mount home0:/disk0/fred /home/fred1 -o fsc,rsize=8192
mount home0:/disk0/fred /home/fred2 -o fsc,rsize=65536
```

- Another way to avoid superblock sharing is to suppress it explicitly with the **nosharecache** parameter. Using the same example:

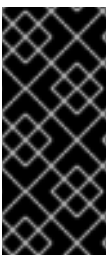
```
mount home0:/disk0/fred /home/fred -o nosharecache,fsc
mount home0:/disk0/jim /home/jim -o nosharecache,fsc
```

However, in this case only one of the superblocks is permitted to use cache since there is nothing to distinguish the Level 2 keys of **home0:/disk0/fred** and **home0:/disk0/jim**.

- To specify the addressing to the superblock, add a *unique identifier* on at least one of the mounts, i.e. **fsc=unique-identifier**.

```
mount home0:/disk0/fred /home/fred -o nosharecache,fsc
mount home0:/disk0/jim /home/jim -o nosharecache,fsc=jim
```

Here, the unique identifier **jim** is added to the Level 2 key used in the cache for **/home/jim**.



IMPORTANT

The user can not share caches between superblocks that have different communications or protocol parameters. For example, it is not possible to share between NFSv4.0 and NFSv3 or between NFSv4.1 and NFSv4.2 because they force different superblocks. Also setting parameters, such as the read size (**rsize**), prevents cache sharing because, again, it forces a different superblock.

10.2. CACHE LIMITATIONS WITH NFS

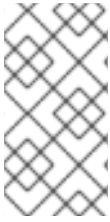
There are some cache limitations with NFS:

- Opening a file from a shared file system for direct I/O automatically bypasses the cache. This is because this type of access must be direct to the server.
- Opening a file from a shared file system for either direct I/O or writing flushes the cached copy of the file. FS-Cache will not cache the file again until it is no longer opened for direct I/O or writing.

- Furthermore, this release of FS-Cache only caches regular NFS files. FS-Cache will *not* cache directories, symlinks, device files, FIFOs and sockets.

CHAPTER 11. MOUNTING AN SMB SHARE ON RED HAT ENTERPRISE LINUX

The Server Message Block (SMB) protocol implements an application-layer network protocol used to access resources on a server, such as file shares and shared printers.



NOTE

In the context of SMB, you can find mentions about the Common Internet File System (CIFS) protocol, which is a dialect of SMB. Both the SMB and CIFS protocol are supported, and the kernel module and utilities involved in mounting SMB and CIFS shares both use the name **cifs**.

This section describes how to mount shares from an SMB server.

Prerequisites

On Microsoft Windows, SMB is implemented by default. On Red Hat Enterprise Linux, the **cifs.ko** file system module of the kernel provides support for mounting SMB shares. Therefore, install the **cifs-utils** package:

```
# dnf install cifs-utils
```

The **cifs-utils** package provides utilities to:

- Mount SMB and CIFS shares
- Manage NT Lan Manager (NTLM) credentials in the kernel's keyring
- Set and display Access Control Lists (ACL) in a security descriptor on SMB and CIFS shares

11.1. SUPPORTED SMB PROTOCOL VERSIONS

The **cifs.ko** kernel module supports the following SMB protocol versions:

- SMB 1



WARNING

The SMB1 protocol is deprecated due to known security issues, and is only **safe to use on a private network**. The main reason that SMB1 is still provided as a supported option is that currently it is the only SMB protocol version that supports UNIX extensions. If you do not need to use UNIX extensions on SMB, Red Hat strongly recommends using SMB2 or later.

- SMB 2.0
- SMB 2.1

- SMB 3.0
- SMB 3.1.1

**NOTE**

Depending on the protocol version, not all SMB features are implemented.

11.2. UNIX EXTENSIONS SUPPORT

Samba uses the **CAP_UNIX** capability bit in the SMB protocol to provide the UNIX extensions feature. These extensions are also supported by the **cifs.ko** kernel module. However, both Samba and the kernel module support UNIX extensions only in the SMB 1 protocol.

To use UNIX extensions:

1. Set the **server min protocol** parameter in the **[global]** section in the **/etc/samba/smb.conf** file to **NT1**.
2. Mount the share using the SMB 1 protocol by providing the **-o vers=1.0** option to the mount command. For example:

```
# mount -t cifs -o vers=1.0,username=user_name //server_name/share_name /mnt/
```

By default, the kernel module uses SMB 2 or the highest later protocol version supported by the server. Passing the **-o vers=1.0** option to the **mount** command forces that the kernel module uses the SMB 1 protocol that is required for using UNIX extensions.

To verify if UNIX extensions are enabled, display the options of the mounted share:

```
# mount
...
//server/share on /mnt type cifs (...,unix,...)
```

If the **unix** entry is displayed in the list of mount options, UNIX extensions are enabled.

11.3. MANUALLY MOUNTING AN SMB SHARE

If you only require an SMB share to be temporary mounted, you can mount it manually using the **mount** utility.

**NOTE**

Manually mounted shares are not mounted automatically again when you reboot the system. To configure that Red Hat Enterprise Linux automatically mounts the share when the system boots, see [Mounting an SMB share automatically when the system boots](#) .

Prerequisites

- The **cifs-utils** package is installed.

Procedure

To manually mount an SMB share, use the **mount** utility with the **-t cifs** parameter:


```
# mount -t cifs -o username=user_name //server_name/share_name /mnt/
Password for user_name@//server_name/share_name: password
```

In the **-o** parameter, you can specify options that are used to mount the share. For details, see the **OPTIONS** section in the **mount.cifs(8)** man page and [Frequently used mount options](#).

Example 11.1. Mounting a share using an encrypted SMB 3.0 connection

To mount the `\\server\example\` share as the **DOMAINAdministrator** user over an encrypted SMB 3.0 connection into the `/mnt/` directory:

```
# mount -t cifs -o username=DOMAINAdministrator,seal,vers=3.0 //server/example /mnt/
Password for DOMAINAdministrator@//server_name/share_name: password
```

11.4. MOUNTING AN SMB SHARE AUTOMATICALLY WHEN THE SYSTEM BOOTS

If access to a mounted SMB share is permanently required on a server, mount the share automatically at boot time.

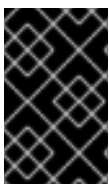
Prerequisites

- The **cifs-utils** package is installed.

Procedure

To mount an SMB share automatically when the system boots, add an entry for the share to the `/etc/fstab` file. For example:

```
//server_name/share_name /mnt cifs credentials=/root/smb.cred 0 0
```



IMPORTANT

To enable the system to mount a share automatically, you must store the user name, password, and domain name in a credentials file. For details, see [Authenticating to an SMB share using a credentials file](#).

In the fourth field of the row in the `/etc/fstab`, specify mount options, such as the path to the credentials file. For details, see the **OPTIONS** section in the **mount.cifs(8)** man page and [Frequently used mount options](#).

To verify that the share mounts successfully, enter:

```
# mount /mnt/
```

11.5. AUTHENTICATING TO AN SMB SHARE USING A CREDENTIALS FILE

In certain situations, such as when mounting a share automatically at boot time, a share should be mounted without entering the user name and password. To implement this, create a credentials file.

Prerequisites

- The **cifs-utils** package is installed.

Procedure

1. Create a file, such as **/root/smb.cred**, and specify the user name, password, and domain name that file:

```
username=user_name
password=password
domain=domain_name
```

2. Set the permissions to only allow the owner to access the file:

```
# chown user_name /root/smb.cred
# chmod 600 /root/smb.cred
```

You can now pass the **credentials=file_name** mount option to the **mount** utility or use it in the **/etc/fstab** file to mount the share without being prompted for the user name and password.

11.6. FREQUENTLY USED MOUNT OPTIONS

When you mount an SMB share, the mount options determine:

- How the connection will be established with the server. For example, which SMB protocol version is used when connecting to the server.
- How the share will be mounted into the local file system. For example, if the system overrides the remote file and directory permissions to enable multiple local users to access the content on the server.

To set multiple options in the fourth field of the **/etc/fstab** file or in the **-o** parameter of a mount command, separate them with commas. For example, see [Mounting a share with the multiuser option](#).

The following list gives frequently used mount options:

Option	Description
<code>credentials=file_name</code>	Sets the path to the credentials file. See Authenticating to an SMB share using a credentials file .
<code>dir_mode=mode</code>	Sets the directory mode if the server does not support CIFS UNIX extensions.
<code>file_mode=mode</code>	Sets the file mode if the server does not support CIFS UNIX extensions.
<code>password=password</code>	Sets the password used to authenticate to the SMB server. Alternatively, specify a credentials file using the credentials option.
<code>seal</code>	Enables encryption support for connections using SMB 3.0 or a later protocol version. Therefore, use seal together with the vers mount option set to 3.0 or later. See the example in Manually mounting an SMB share .

Option	Description
<code>sec=security_mode</code>	<p>Sets the security mode, such as ntlmsspi, to enable NTLMv2 password hashing and enabled packet signing. For a list of supported values, see the option's description in the mount.cifs(8) man page.</p> <p>If the server does not support the ntlmv2 security mode, use sec=ntlmssp, which is the default.</p> <p>For security reasons, do not use the insecure ntlm security mode.</p>
<code>username=user_name</code>	<p>Sets the user name used to authenticate to the SMB server. Alternatively, specify a credentials file using the credentials option.</p>
<code>vers=SMB_protocol_version</code>	<p>Sets the SMB protocol version used for the communication with the server.</p>

For a complete list, see the **OPTIONS** section in the **mount.cifs(8)** man page.

CHAPTER 12. PERFORMING A MULTI-USER SMB MOUNT

The credentials you provide to mount a share determine the access permissions on the mount point by default. For example, if you use the **DOMAINexample** user when you mount a share, all operations on the share will be executed as this user, regardless which local user performs the operation.

However, in certain situations, the administrator wants to mount a share automatically when the system boots, but users should perform actions on the share's content using their own credentials. The **multiuser** mount options lets you configure this scenario.



IMPORTANT

To use the **multiuser** mount option, you must additionally set the **sec** mount option to a security type that supports providing credentials in a non-interactive way, such as **krb5** or the **ntlmssp** option with a credentials file. For details, see [Accessing a share as a user](#).

The **root** user mounts the share using the **multiuser** option and an account that has minimal access to the contents of the share. Regular users can then provide their user name and password to the current session's kernel keyring using the **cifscreds** utility. If the user accesses the content of the mounted share, the kernel uses the credentials from the kernel keyring instead of the one initially used to mount the share.

Using this feature consists of the following steps:

- [Mount a share with the **multiuser** option.](#)
- [Optionally, verify if the share was successfully mounted with the **multiuser** option.](#)
- [Access the share as a user](#).

Prerequisites

- The **cifs-utils** package is installed.

12.1. MOUNTING A SHARE WITH THE MULTIUSER OPTION

Before users can access the share with their own credentials, mount the share as the **root** user using an account with limited permissions.

Procedure

To mount a share automatically with the **multiuser** option when the system boots:

1. Create the entry for the share in the **/etc/fstab** file. For example:

```
//server_name/share_name /mnt cifs multiuser,sec=ntlmssp,credentials=/root/smb.cred
0 0
```

2. Mount the share:

```
# mount /mnt/
```

If you do not want to mount the share automatically when the system boots, mount it manually by passing **-o multiuser,sec=security_type** to the **mount** command. For details about mounting an SMB share manually, see [Manually mounting an SMB share](#).

12.2. VERIFYING IF AN SMB SHARE IS MOUNTED WITH THE MULTIUSER OPTION

To verify if a share is mounted with the **multiuser** option, display the mount options.

Procedure

```
# mount
...
//server_name/share_name on /mnt type cifs (sec=ntlmssp,multiuser,...)
```

If the **multiuser** entry is displayed in the list of mount options, the feature is enabled.

12.3. ACCESSING A SHARE AS A USER

If an SMB share is mounted with the **multiuser** option, users can provide their credentials for the server to the kernel's keyring:

```
# cifscreds add -u SMB_user_name server_name
Password: password
```

When the user performs operations in the directory that contains the mounted SMB share, the server applies the file system permissions for this user, instead of the one initially used when the share was mounted.



NOTE

Multiple users can perform operations using their own credentials on the mounted share at the same time.

CHAPTER 13. OVERVIEW OF PERSISTENT NAMING ATTRIBUTES

As a system administrator, you need to refer to storage volumes using persistent naming attributes to build storage setups that are reliable over multiple system boots.

13.1. DISADVANTAGES OF NON-PERSISTENT NAMING ATTRIBUTES

Red Hat Enterprise Linux provides a number of ways to identify storage devices. It is important to use the correct option to identify each device when used in order to avoid inadvertently accessing the wrong device, particularly when installing to or reformatting drives.

Traditionally, non-persistent names in the form of `/dev/sd(major number)(minor number)` are used on Linux to refer to storage devices. The major and minor number range and associated **sd** names are allocated for each device when it is detected. This means that the association between the major and minor number range and associated **sd** names can change if the order of device detection changes.

Such a change in the ordering might occur in the following situations:

- The parallelization of the system boot process detects storage devices in a different order with each system boot.
- A disk fails to power up or respond to the SCSI controller. This results in it not being detected by the normal device probe. The disk is not accessible to the system and subsequent devices will have their major and minor number range, including the associated **sd** names shifted down. For example, if a disk normally referred to as **sdb** is not detected, a disk that is normally referred to as **sdc** would instead appear as **sdb**.
- A SCSI controller (host bus adapter, or HBA) fails to initialize, causing all disks connected to that HBA to not be detected. Any disks connected to subsequently probed HBAs are assigned different major and minor number ranges, and different associated **sd** names.
- The order of driver initialization changes if different types of HBAs are present in the system. This causes the disks connected to those HBAs to be detected in a different order. This might also occur if HBAs are moved to different PCI slots on the system.
- Disks connected to the system with Fibre Channel, iSCSI, or FCoE adapters might be inaccessible at the time the storage devices are probed, due to a storage array or intervening switch being powered off, for example. This might occur when a system reboots after a power failure, if the storage array takes longer to come online than the system take to boot. Although some Fibre Channel drivers support a mechanism to specify a persistent SCSI target ID to WWPN mapping, this does not cause the major and minor number ranges, and the associated **sd** names to be reserved; it only provides consistent SCSI target ID numbers.

These reasons make it undesirable to use the major and minor number range or the associated **sd** names when referring to devices, such as in the `/etc/fstab` file. There is the possibility that the wrong device will be mounted and data corruption might result.

Occasionally, however, it is still necessary to refer to the **sd** names even when another mechanism is used, such as when errors are reported by a device. This is because the Linux kernel uses **sd** names (and also SCSI host/channel/target/LUN tuples) in kernel messages regarding the device.

13.2. FILE SYSTEM AND DEVICE IDENTIFIERS

This sections explains the difference between persistent attributes identifying file systems and block devices.

File system identifiers

File system identifiers are tied to a particular file system created on a block device. The identifier is also stored as part of the file system. If you copy the file system to a different device, it still carries the same file system identifier. On the other hand, if you rewrite the device, such as by formatting it with the **mkfs** utility, the device loses the attribute.

File system identifiers include:

- Unique identifier (UUID)
- Label

Device identifiers

Device identifiers are tied to a block device: for example, a disk or a partition. If you rewrite the device, such as by formatting it with the **mkfs** utility, the device keeps the attribute, because it is not stored in the file system.

Device identifiers include:

- World Wide Identifier (WWID)
- Partition UUID
- Serial number

Recommendations

- Some file systems, such as logical volumes, span multiple devices. Red Hat recommends accessing these file systems using file system identifiers rather than device identifiers.

13.3. DEVICE NAMES MANAGED BY THE UDEV MECHANISM IN /DEV/DISK/

The **udev** mechanism is used for all types of devices in Linux, and is not limited only for storage devices. It provides different kinds of persistent naming attributes in the **/dev/disk/** directory. In the case of storage devices, Red Hat Enterprise Linux contains **udev** rules that create symbolic links in the **/dev/disk/** directory. This enables you to refer to storage devices by:

- Their content
- A unique identifier
- Their serial number.

Although **udev** naming attributes are persistent, in that they do not change on their own across system reboots, some are also configurable.

13.3.1. File system identifiers

The UUID attribute in **/dev/disk/by-uuid/**

Entries in this directory provide a symbolic name that refers to the storage device by a **unique identifier** (UUID) in the content (that is, the data) stored on the device. For example:

`/dev/disk/by-uuid/3e6be9de-8139-11d1-9106-a43f08d823a6`

You can use the UUID to refer to the device in the `/etc/fstab` file using the following syntax:

`UUID=3e6be9de-8139-11d1-9106-a43f08d823a6`

You can configure the UUID attribute when creating a file system, and you can also change it later on.

The Label attribute in `/dev/disk/by-label/`

Entries in this directory provide a symbolic name that refers to the storage device by a **label** in the content (that is, the data) stored on the device.

For example:

`/dev/disk/by-label/Boot`

You can use the label to refer to the device in the `/etc/fstab` file using the following syntax:

`LABEL=Boot`

You can configure the Label attribute when creating a file system, and you can also change it later on.

13.3.2. Device identifiers

The WWID attribute in `/dev/disk/by-id/`

The World Wide Identifier (WWID) is a persistent, **system-independent identifier** 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. The identifier is a property of the device but is not stored in the content (that is, the data) on the devices.

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**).

Red Hat Enterprise Linux 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.



NOTE

If your are using an NVMe device, you might run into a disk by-id naming change for some vendors, if the serial number of your device has leading whitespace.

Example 13.1. WWID mappings

WWID symlink	Non-persistent device	Note
<code>/dev/disk/by-id/scsi-3600508b400105e210000900000490000</code>	<code>/dev/sda</code>	A device with a page 0x83 identifier

WWID symlink	Non-persistent device	Note
/dev/disk/by-id/scsi-SSEAGATE_ST373453LW_3HW1RHM6	/dev/sdb	A device with a page 0x80 identifier
/dev/disk/by-id/ata-SAMSUNG_MZNLN256MHQ-000L7_S2WDNX0J336519-part3	/dev/sdc3	A disk partition

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.

The Partition UUID attribute in /dev/disk/by-partuuid

The Partition UUID (PARTUUID) attribute identifies partitions as defined by GPT partition table.

Example 13.2. Partition UUID mappings

PARTUUID symlink	Non-persistent device
/dev/disk/by-partuuid/4cd1448a-01	/dev/sda1
/dev/disk/by-partuuid/4cd1448a-02	/dev/sda2
/dev/disk/by-partuuid/4cd1448a-03	/dev/sda3

The Path attribute in /dev/disk/by-path/

This attribute provides a symbolic name that refers to the storage device by the **hardware path** used to access the device.

The Path attribute fails if any part of the hardware path (for example, the PCI ID, target port, or LUN number) changes. The Path attribute is therefore unreliable. However, the Path attribute may be useful in one of the following scenarios:

- You need to identify a disk that you are planning to replace later.
- You plan to install a storage service on a disk in a specific location.

13.4. THE WORLD WIDE IDENTIFIER WITH DM MULTIPATH

You can configure Device Mapper (DM) Multipath to map between the World Wide Identifier (WWID) and non-persistent device names.

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

The command **multipath -l** shows the mapping to the non-persistent identifiers:

- **Host:Channel:Target:LUN**
- **/dev/sd** name
- **major:minor** number

Example 13.3. WWID mappings in a multipath configuration

An example output of the **multipath -l** command:

```
3600508b400105df70000e00000ac0000 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]
```

DM 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 DM 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 **/etc/multipath/bindings**. These **mpathN** names are persistent as long as that file is maintained.



IMPORTANT

If you use **user_friendly_names**, then additional steps are required to obtain consistent names in a cluster.

13.5. LIMITATIONS OF THE UDEV DEVICE NAMING CONVENTION

The following are some limitations of the **udev** naming convention:

- It is possible that the device might not be accessible at the time the query is performed because the **udev** mechanism might rely on the ability to query the storage device when the **udev** rules are processed for a **udev** event. This is more likely to occur with Fibre Channel, iSCSI or FCoE storage devices when the device is not located in the server chassis.
- The kernel might send **udev** events at any time, causing the rules to be processed and possibly causing the **/dev/disk/by-*/** links to be removed if the device is not accessible.
- There might be a delay between when the **udev** event is generated and when it is processed, such as when a large number of devices are detected and the user-space **udev** service takes some amount of time to process the rules for each one. This might cause a delay between when the kernel detects the device and when the **/dev/disk/by-*/** names are available.
- External programs such as **blkid** invoked by the rules might open the device for a brief period of time, making the device inaccessible for other uses.
- The device names managed by the **udev** mechanism in **/dev/disk/** may change between major releases, requiring you to update the links.

13.6. LISTING PERSISTENT NAMING ATTRIBUTES

This procedure describes how to find out the persistent naming attributes of non-persistent storage devices.

Procedure

- To list the UUID and Label attributes, use the **lsblk** utility:

```
$ lsblk --fs storage-device
```

For example:

Example 13.4. Viewing the UUID and Label of a file system

```
$ lsblk --fs /dev/sda1
```

NAME	FSTYPE	LABEL	UUID	MOUNTPOINT
sda1	xf	Boot	afa5d5e3-9050-48c3-acc1-bb30095f3dc4	/boot

- To list the PARTUUID attribute, use the **lsblk** utility with the **--output +PARTUUID** option:

```
$ lsblk --output +PARTUUID
```

For example:

Example 13.5. Viewing the PARTUUID attribute of a partition

```
$ lsblk --output +PARTUUID /dev/sda1
```

NAME	MAJ:MIN	RM	SIZE	RO	TYPE	MOUNTPOINT	PARTUUID
sda1	8:1	0	512M	0	part	/boot	4cd1448a-01

- To list the WWID attribute, examine the targets of symbolic links in the **/dev/disk/by-id/** directory. For example:

Example 13.6. Viewing the WWID of all storage devices on the system

```
$ file /dev/disk/by-id/*
```

```
/dev/disk/by-id/ata-QEMU_HARDDISK_QM00001
symbolic link to ../../sda
/dev/disk/by-id/ata-QEMU_HARDDISK_QM00001-part1
symbolic link to ../../sda1
/dev/disk/by-id/ata-QEMU_HARDDISK_QM00001-part2
symbolic link to ../../sda2
/dev/disk/by-id/dm-name-rhel_rhel8-root
symbolic link to ../../dm-0
/dev/disk/by-id/dm-name-rhel_rhel8-swap
symbolic link to ../../dm-1
/dev/disk/by-id/dm-uuid-LVM-
QIWtEHtXGobe5bewlIUDivKOz5ofkgFhP0RMFsNyySVihqEl2cWWbR7MjXJolD6g
```

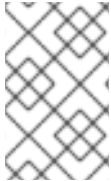
```

symbolic link to ../../dm-1
/dev/disk/by-id/dm-uuid-LVM-
QIWtEHtXGobe5bewllUDivKOz5ofkgFhXqH2M45hD2H9nAf2qfWSrlRLhzfMyOKd
symbolic link to ../../dm-0
/dev/disk/by-id/lvm-pv-uuid-atlr2Y-vuMo-ueoH-CpMG-4JuH-AhEF-wu4QQm
symbolic link to ../../sda2

```

13.7. MODIFYING PERSISTENT NAMING ATTRIBUTES

This procedure describes how to change the UUID or Label persistent naming attribute of a file system.



NOTE

Changing **udev** attributes happens in the background and might take a long time. The **udevadm settle** command waits until the change is fully registered, which ensures that your next command will be able to utilize the new attribute correctly.

In the following commands:

- Replace *new-uuid* with the UUID you want to set; for example, **1cdfbc07-1c90-4984-b5ec-f61943f5ea50**. You can generate a UUID using the **uuidgen** command.
- Replace *new-label* with a label; for example, **backup_data**.

Prerequisites

- If you are modifying the attributes of an XFS file system, unmount it first.

Procedure

- To change the UUID or Label attributes of an **XFS** file system, use the **xfs_admin** utility:

```

# xfs_admin -U new-uuid -L new-label storage-device
# udevadm settle

```

- To change the UUID or Label attributes of an **ext4**, **ext3**, or **ext2** file system, use the **tune2fs** utility:

```

# tune2fs -U new-uuid -L new-label storage-device
# udevadm settle

```

- To change the UUID or Label attributes of a swap volume, use the **swaponlabel** utility:

```

# swaponlabel --uuid new-uuid --label new-label swap-device
# udevadm settle

```

CHAPTER 14. PARTITION OPERATIONS WITH PARTED

parted is a program to manipulate disk partitions. It supports multiple partition table formats, including MS-DOS and GPT. It is useful for creating space for new operating systems, reorganizing disk usage, and copying data to new hard disks.

14.1. VIEWING THE PARTITION TABLE WITH PARTED

Display the partition table of a block device to see the partition layout and details about individual partitions. You can view the partition table on a block device using the **parted** utility.

Procedure

1. Start the **parted** utility. For example, the following output lists the device **/dev/sda**:

```
# parted /dev/sda
```

2. View the partition table:

```
# (parted) print

Model: ATA SAMSUNG MZNLN256 (scsi)
Disk /dev/sda: 256GB
Sector size (logical/physical): 512B/512B
Partition Table: msdos
Disk Flags:

Number Start End   Size Type  File system  Flags
 1    1049kB 269MB 268MB primary xfs          boot
 2    269MB 34.6GB 34.4GB primary
 3    34.6GB 45.4GB 10.7GB primary
 4    45.4GB 256GB 211GB extended
 5    45.4GB 256GB 211GB logical
```

3. Optional: Switch to the device you want to examine next:

```
# (parted) select block-device
```

For a detailed description of the print command output, see the following:

Model: ATA SAMSUNG MZNLN256 (scsi)

The disk type, manufacturer, model number, and interface.

Disk /dev/sda: 256GB

The file path to the block device and the storage capacity.

Partition Table: msdos

The disk label type.

Number

The partition number. For example, the partition with minor number 1 corresponds to **/dev/sda1**.

Start and End

The location on the device where the partition starts and ends.

Type

Valid types are metadata, free, primary, extended, or logical.

File system

The file system type. If the **File system** field of a device shows no value, this means that its file system type is unknown. The **parted** utility cannot recognize the file system on encrypted devices.

Flags

Lists the flags set for the partition. Available flags are **boot**, **root**, **swap**, **hidden**, **raid**, **lvm**, or **lba**.

Additional resources

- **parted(8)** man page.

14.2. CREATING A PARTITION TABLE ON A DISK WITH PARTED

Use the **parted** utility to format a block device with a partition table more easily.



WARNING

Formatting a block device with a partition table deletes all data stored on the device.

Procedure

1. Start the interactive **parted** shell:

```
# parted block-device
```

2. Determine if there already is a partition table on the device:

```
# (parted) print
```

If the device already contains partitions, they will be deleted in the following steps.

3. Create the new partition table:

```
# (parted) mklabel table-type
```

- Replace *table-type* with the intended partition table type:
 - **msdos** for MBR
 - **gpt** for GPT

Example 14.1. Creating a GUID Partition Table (GPT) table

To create a GPT table on the disk, use:

```
# (parted) mklabel gpt
```

The changes start applying after you enter this command.

4. View the partition table to confirm that it is created:

```
# (parted) print
```

5. Exit the **parted** shell:

```
# (parted) quit
```

Additional resources

- **parted(8)** man page.

14.3. CREATING A PARTITION WITH PARTED

As a system administrator, you can create new partitions on a disk by using the **parted** utility.



NOTE

The required partitions are **swap**, **/boot/**, and **/ (root)**.

Prerequisites

- A partition table on the disk.
- If the partition you want to create is larger than 2TiB, format the disk with the **GUID Partition Table (GPT)**.

Procedure

1. Start the **parted** utility:

```
# parted block-device
```

2. View the current partition table to determine if there is enough free space:

```
# (parted) print
```

- Resize the partition in case there is not enough free space.
 - From the partition table, determine:
 - The start and end points of the new partition.
 - On MBR, what partition type it should be.
3. Create the new partition:

```
# (parted) mkpart part-type name fs-type start end
```

- Replace *part-type* with **primary**, **logical**, or **extended**. This applies only to the MBR partition table.
- Replace *name* with an arbitrary partition name. This is required for GPT partition tables.
- Replace *fs-type* with **xfs**, **ext2**, **ext3**, **ext4**, **fat16**, **fat32**, **hfs**, **hfs+**, **linux-swaps**, **ntfs**, or **reiserfs**. The *fs-type* parameter is optional. Note that the **parted** utility does not create the file system on the partition.
- Replace *start* and *end* with the sizes that determine the starting and ending points of the partition, counting from the beginning of the disk. You can use size suffixes, such as **512MiB**, **20GiB**, or **1.5TiB**. The default size is in megabytes.

Example 14.2. Creating a small primary partition

To create a primary partition from 1024MiB until 2048MiB on an MBR table, use:

```
# (parted) mkpart primary 1024MiB 2048MiB
```

The changes start applying after you enter the command.

4. View the partition table to confirm that the created partition is in the partition table with the correct partition type, file system type, and size:

```
# (parted) print
```

5. Exit the **parted** shell:

```
# (parted) quit
```

6. Register the new device node:

```
# udevadm settle
```

7. Verify that the kernel recognizes the new partition:

```
# cat /proc/partitions
```

Additional resources

- **parted(8)** man page.
- [Creating a partition table on a disk with parted](#) .
- [Resizing a partition with parted](#)

14.4. REMOVING A PARTITION WITH PARTED

Using the **parted** utility, you can remove a disk partition to free up disk space.

**WARNING**

Removing a partition deletes all data stored on the partition.

Procedure

1. Start the interactive **parted** shell:

```
# parted block-device
```

- Replace *block-device* with the path to the device where you want to remove a partition: for example, **/dev/sda**.

2. View the current partition table to determine the minor number of the partition to remove:

```
(parted) print
```

3. Remove the partition:

```
(parted) rm minor-number
```

- Replace *minor-number* with the minor number of the partition you want to remove.

The changes start applying as soon as you enter this command.

4. Verify that you have removed the partition from the partition table:

```
(parted) print
```

5. Exit the **parted** shell:

```
(parted) quit
```

6. Verify that the kernel registers that the partition is removed:

```
# cat /proc/partitions
```

7. Remove the partition from the **/etc/fstab** file, if it is present. Find the line that declares the removed partition, and remove it from the file.

8. Regenerate mount units so that your system registers the new **/etc/fstab** configuration:

```
# systemctl daemon-reload
```

9. If you have deleted a swap partition or removed pieces of LVM, remove all references to the partition from the kernel command line:

- a. List active kernel options and see if any option references the removed partition:

```
# grubby --info=ALL
```

- b. Remove the kernel options that reference the removed partition:

```
# grubby --update-kernel=ALL --remove-args="option"
```

10. To register the changes in the early boot system, rebuild the **initramfs** file system:

```
# dracut --force --verbose
```

Additional resources

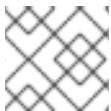
- **parted(8)** man page

14.5. RESIZING A PARTITION WITH PARTED

Using the **parted** utility, extend a partition to use unused disk space, or shrink a partition to use its capacity for different purposes.

Prerequisites

- Back up the data before shrinking a partition.
- If the partition you want to create is larger than 2TiB, format the disk with the **GUID Partition Table (GPT)**.
- If you want to shrink the partition, first shrink the file system so that it is not larger than the resized partition.



NOTE

XFS does not support shrinking.

Procedure

1. Start the **parted** utility:

```
# parted block-device
```

2. View the current partition table:

```
# (parted) print
```

From the partition table, determine:

- The minor number of the partition.
 - The location of the existing partition and its new ending point after resizing.
3. Resize the partition:

```
# (parted) resizepart 1 2GiB
```

- Replace 1 with the minor number of the partition that you are resizing.
 - Replace 2 with the size that determines the new ending point of the resized partition, counting from the beginning of the disk. You can use size suffixes, such as **512MiB**, **20GiB**, or **1.5TiB**. The default size is in megabytes.
4. View the partition table to confirm that the resized partition is in the partition table with the correct size:

```
# (parted) print
```

5. Exit the **parted** shell:

```
# (parted) quit
```

6. Verify that the kernel registers the new partition:

```
# cat /proc/partitions
```

7. Optional: If you extended the partition, extend the file system on it as well.

Additional resources

- **parted(8)** man page.
- [Creating a partition table on a disk with parted](#)
- [Resizing an ext4 file system](#)
- [Increasing the size of an XFS file system](#)

CHAPTER 15. STRATEGIES FOR REPARTITIONING A DISK

There are different approaches to repartitioning a disk. These include:

- Unpartitioned free space is available.
- An unused partition is available.
- Free space in an actively used partition is available.



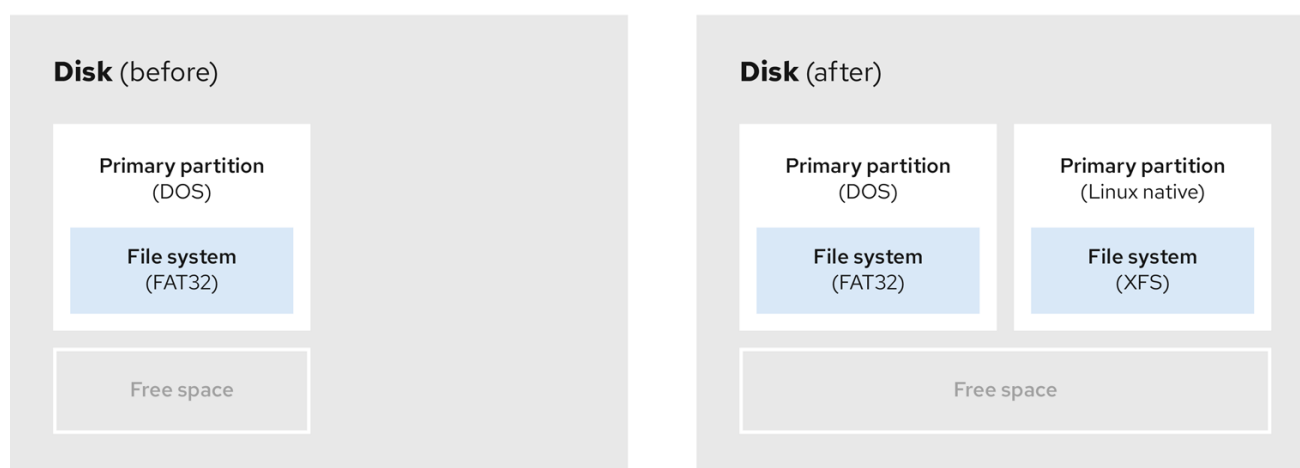
NOTE

The following examples are simplified for clarity and do not reflect the exact partition layout when actually installing Red Hat Enterprise Linux.

15.1. USING UNPARTITIONED FREE SPACE

Partitions that are already defined and do not span the entire hard disk, leave unallocated space that is not part of any defined partition. The following diagram shows what this might look like.

Figure 15.1. Disk with unpartitioned free space



269_RHEL_0822

The first diagram represents a disk with one primary partition and an undefined partition with unallocated space. The second diagram represents a disk with two defined partitions with allocated space.

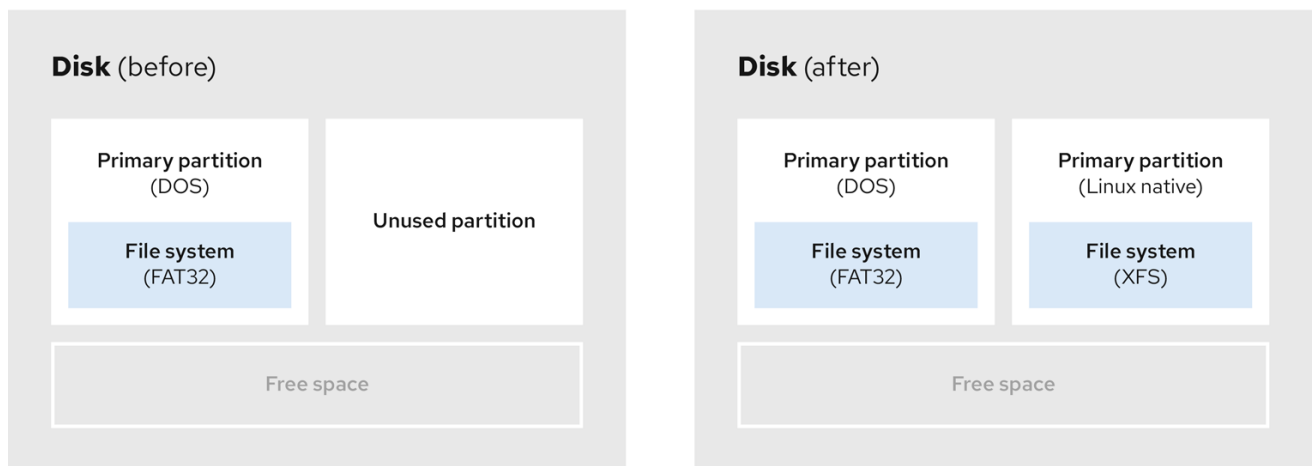
An unused hard disk also falls into this category. The only difference is that *all* the space is not part of any defined partition.

On a new disk, you can create the necessary partitions from the unused space. Most preinstalled operating systems are configured to take up all available space on a disk drive.

15.2. USING SPACE FROM AN UNUSED PARTITION

In the following example, the first diagram represents a disk with an unused partition. The second diagram represents reallocating an unused partition for Linux.

Figure 15.2. Disk with an unused partition



269_RHEL_0822

To use the space allocated to the unused partition, delete the partition and then create the appropriate Linux partition instead. Alternatively, during the installation process, delete the unused partition and manually create new partitions.

15.3. USING FREE SPACE FROM AN ACTIVE PARTITION

This process can be difficult to manage because an active partition, that is already in use, contains the required free space. In most cases, hard disks of computers with preinstalled software contain one larger partition holding the operating system and data.



WARNING

If you want to use an operating system (OS) on an active partition, you must reinstall the OS. Be aware that some computers, which include pre-installed software, do not include installation media to reinstall the original OS. Check whether this applies to your OS before you destroy an original partition and the OS installation.

To optimise the use of available free space, you can use the methods of destructive or non-destructive repartitioning.

15.3.1. Destructive repartitioning

Destructive repartitioning destroys the partition on your hard drive and creates several smaller partitions instead. Backup any needed data from the original partition as this method deletes the complete contents.

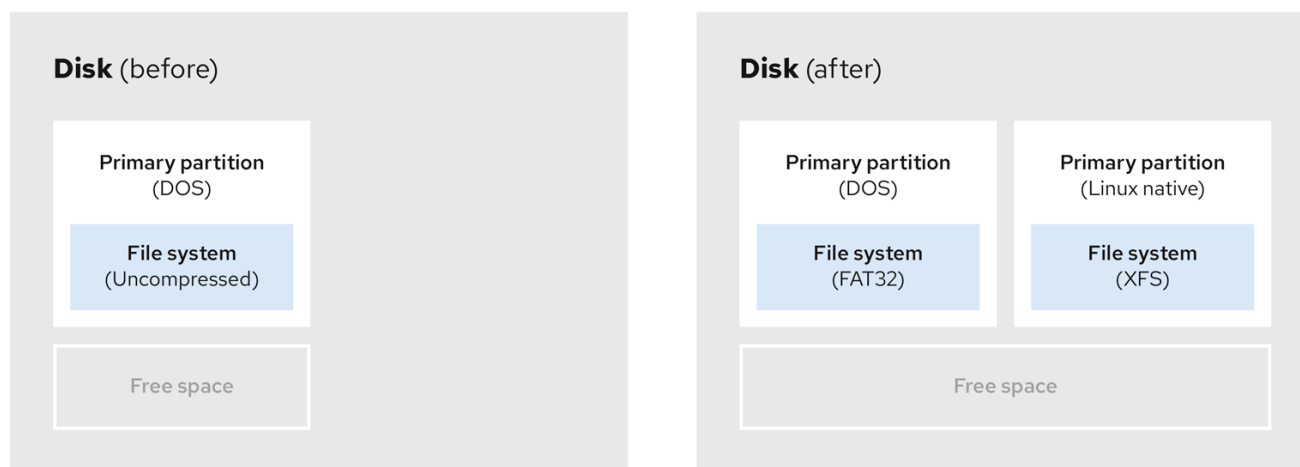
After creating a smaller partition for your existing operating system, you can:

- Reinstall software.
- Restore your data.

- Start your Red Hat Enterprise Linux installation.

The following diagram is a simplified representation of using the destructive repartitioning method.

Figure 15.3. Destructive repartitioning action on disk



269_RHEL_0822



WARNING

This method deletes all data previously stored in the original partition.

15.3.2. Non-destructive repartitioning

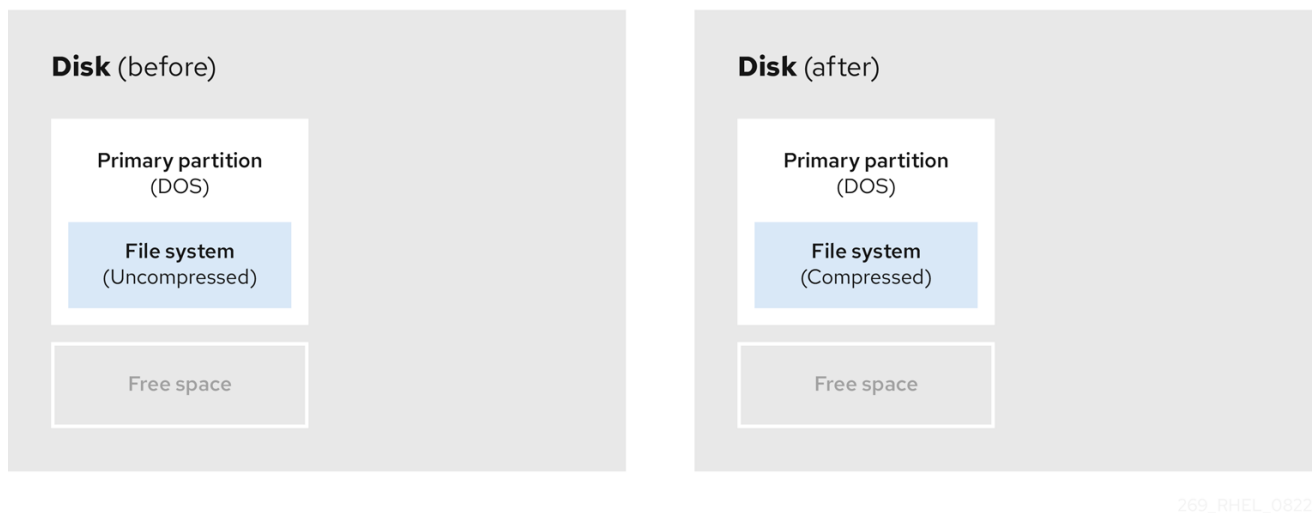
Non-destructive repartitioning resizes partitions, without any data loss. This method is reliable, however it takes longer processing time on large drives.

The following is a list of methods, which can help initiate non-destructive repartitioning.

- Compress existing data

The storage location of some data cannot be changed. This can prevent the resizing of a partition to the required size, and ultimately lead to a destructive repartition process. Compressing data in an already existing partition can help you resize your partitions as needed. It can also help to maximize the free space available.

The following diagram is a simplified representation of this process.

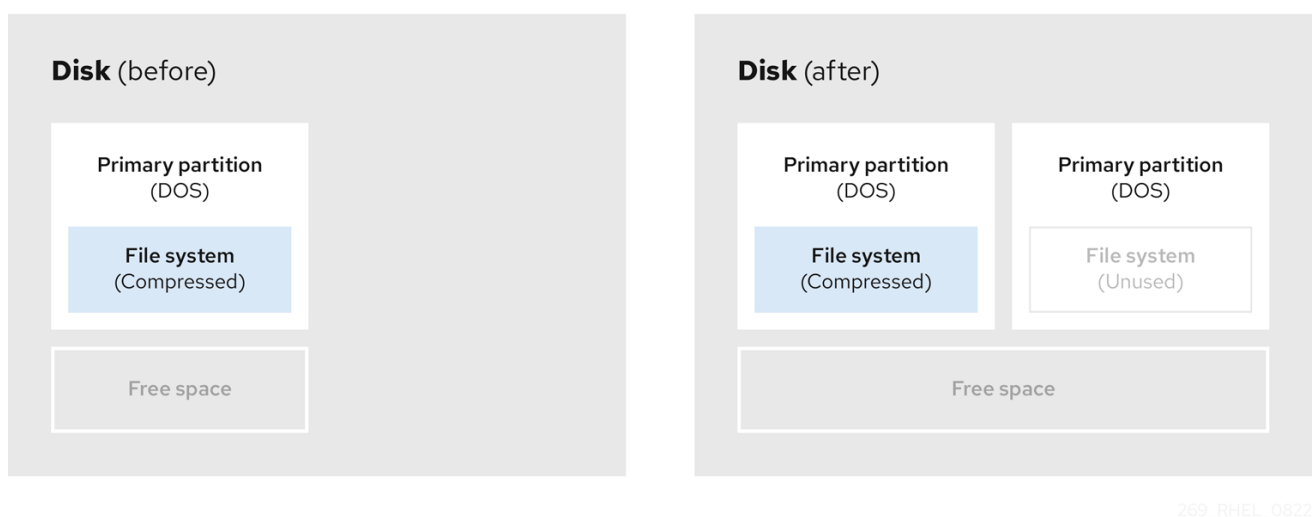
Figure 15.4. Data compression on a disk

To avoid any possible data loss, create a backup before continuing with the compression process.

- Resize the existing partition

By resizing an already existing partition, you can free up more space. Depending on your resizing software, the results may vary. In the majority of cases, you can create a new unformatted partition of the same type, as the original partition.

The steps you take after resizing can depend on the software you use. In the following example, the best practice is to delete the new DOS (Disk Operating System) partition, and create a Linux partition instead. Verify what is most suitable for your disk before initiating the resizing process.

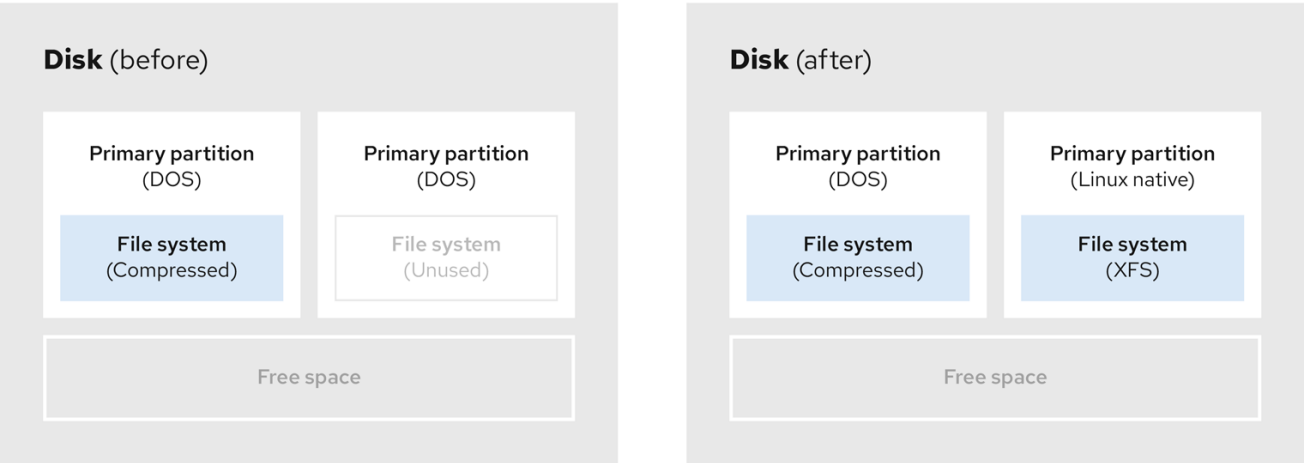
Figure 15.5. Partition resizing on a disk

- Optional: Create new partitions

Some pieces of resizing software support Linux based systems. In such cases, there is no need to delete the newly created partition after resizing. Creating a new partition afterwards depends on the software you use.

The following diagram represents the disk state, before and after creating a new partition.

Figure 15.6. Disk with final partition configuration



269_RHEL_0822

CHAPTER 16. GETTING STARTED WITH XFS

This is an overview of how to create and maintain XFS file systems.

16.1. THE XFS FILE SYSTEM

XFS is a highly scalable, high-performance, robust, and mature 64-bit journaling file system that supports very large files and file systems on a single host. It is the default file system in Red Hat Enterprise Linux 9. XFS was originally developed in the early 1990s by SGI and has a long history of running on extremely large servers and storage arrays.

The features of XFS include:

Reliability

- Metadata journaling, which ensures file system integrity after a system crash by keeping a record of file system operations that can be replayed when the system is restarted and the file system remounted
- Extensive run-time metadata consistency checking
- Scalable and fast repair utilities
- Quota journaling. This avoids the need for lengthy quota consistency checks after a crash.

Scalability and performance

- Supported file system size up to 1024 TiB
- Ability to support a large number of concurrent operations
- B-tree indexing for scalability of free space management
- Sophisticated metadata read-ahead algorithms
- Optimizations for streaming video workloads

Allocation schemes

- Extent-based allocation
- Stripe-aware allocation policies
- Delayed allocation
- Space pre-allocation
- Dynamically allocated inodes

Other features

- Reflink-based file copies
- Tightly integrated backup and restore utilities
- Online defragmentation

- Online file system growing
- Comprehensive diagnostics capabilities
- Extended attributes (**xattr**). This allows the system to associate several additional name/value pairs per file.
- Project or directory quotas. This allows quota restrictions over a directory tree.
- Subsecond timestamps

Performance characteristics

XFS has a high performance on large systems with enterprise workloads. A large system is one with a relatively high number of CPUs, multiple HBAs, and connections to external disk arrays. XFS also performs well on smaller systems that have a multi-threaded, parallel I/O workload.

XFS has a relatively low performance for single threaded, metadata-intensive workloads: for example, a workload that creates or deletes large numbers of small files in a single thread.

16.2. COMPARISON OF TOOLS USED WITH EXT4 AND XFS

This section compares which tools to use to accomplish common tasks on the ext4 and XFS file systems.

Task	ext4	XFS
Create a file system	mkfs.ext4	mkfs.xfs
File system check	e2fsck	xfs_repair
Resize a file system	resize2fs	xfs_growfs
Save an image of a file system	e2image	xfs_metadump and xfs_mdrestore
Label or tune a file system	tune2fs	xfs_admin
Back up a file system	tar and rsync	xfsdump and xfsrestore
Quota management	quota	xfs_quota
File mapping	filefrag	xfs_bmap



NOTE

If you want a complete client-server solution for backups over network, you can use **bacula** backup utility that is available in RHEL 9. For more information about Bacula, see [Bacula backup solution](#).

CHAPTER 17. CREATING AN XFS FILE SYSTEM

As a system administrator, you can create an XFS file system on a block device to enable it to store files and directories.

17.1. CREATING AN XFS FILE SYSTEM WITH MKFS.XFS

This procedure describes how to create an XFS file system on a block device.

Procedure

1. To create the file system:

- If the device is a regular partition, an LVM volume, an MD volume, a disk, or a similar device, use the following command:

```
# mkfs.xfs block-device
```

- Replace *block-device* with the path to the block device. For example, `/dev/sdb1`, `/dev/disk/by-uuid/05e99ec8-def1-4a5e-8a9d-5945339ceb2a`, or `/dev/my-volgroup/my-lv`.
- In general, the default options are optimal for common use.
- When using **mkfs.xfs** on a block device containing an existing file system, add the **-f** option to overwrite that file system.
- To create the file system on a hardware RAID device, check if the system correctly detects the stripe geometry of the device:
 - If the stripe geometry information is correct, no additional options are needed. Create the file system:

```
# mkfs.xfs block-device
```

- If the information is incorrect, specify stripe geometry manually with the **su** and **sw** parameters of the **-d** option. The **su** parameter specifies the RAID chunk size, and the **sw** parameter specifies the number of data disks in the RAID device.

For example:

```
# mkfs.xfs -d su=64k,sw=4 /dev/sda3
```

2. Use the following command to wait for the system to register the new device node:

```
# udevadm settle
```

Additional resources

- **mkfs.xfs(8)** man page.

CHAPTER 18. BACKING UP AN XFS FILE SYSTEM

As a system administrator, you can use the **xfsdump** to back up an XFS file system into a file or on a tape. This provides a simple backup mechanism.

18.1. FEATURES OF XFS BACKUP

This section describes key concepts and features of backing up an XFS file system with the **xfsdump** utility.

You can use the **xfsdump** utility to:

- Perform backups to regular file images.
Only one backup can be written to a regular file.
- Perform backups to tape drives.
The **xfsdump** utility also enables you to write multiple backups to the same tape. A backup can span multiple tapes.

To back up multiple file systems to a single tape device, simply write the backup to a tape that already contains an XFS backup. This appends the new backup to the previous one. By default, **xfsdump** never overwrites existing backups.

- Create incremental backups.
The **xfsdump** utility uses dump levels to determine a base backup to which other backups are relative. Numbers from 0 to 9 refer to increasing dump levels. An incremental backup only backs up files that have changed since the last dump of a lower level:
 - To perform a full backup, perform a level 0 dump on the file system.
 - A level 1 dump is the first incremental backup after a full backup. The next incremental backup would be level 2, which only backs up files that have changed since the last level 1 dump; and so on, to a maximum of level 9.
- Exclude files from a backup using size, subtree, or inode flags to filter them.

Additional resources

- **xfsdump(8)** man page.

18.2. BACKING UP AN XFS FILE SYSTEM WITH XFSDUMP

This procedure describes how to back up the content of an XFS file system into a file or a tape.

Prerequisites

- An XFS file system that you can back up.
- Another file system or a tape drive where you can store the backup.

Procedure

- Use the following command to back up an XFS file system:

```
# xfsdump -l level [-L label] \
    -f backup-destination path-to-xfs-filesystem
```

- Replace *level* with the dump level of your backup. Use **0** to perform a full backup or **1** to **9** to perform consequent incremental backups.
- Replace *backup-destination* with the path where you want to store your backup. The destination can be a regular file, a tape drive, or a remote tape device. For example, **/backup-files/Data.xfsdump** for a file or **/dev/st0** for a tape drive.
- Replace *path-to-xfs-filesystem* with the mount point of the XFS file system you want to back up. For example, **/mnt/data/**. The file system must be mounted.
- When backing up multiple file systems and saving them on a single tape device, add a session label to each backup using the **-L *label*** option so that it is easier to identify them when restoring. Replace *label* with any name for your backup: for example, **backup_data**.

Example 18.1. Backing up multiple XFS file systems

- To back up the content of XFS file systems mounted on the **/boot/** and **/data/** directories and save them as files in the **/backup-files/** directory:

```
# xfsdump -l 0 -f /backup-files/boot.xfsdump /boot
# xfsdump -l 0 -f /backup-files/data.xfsdump /data
```

- To back up multiple file systems on a single tape device, add a session label to each backup using the **-L *label*** option:

```
# xfsdump -l 0 -L "backup_boot" -f /dev/st0 /boot
# xfsdump -l 0 -L "backup_data" -f /dev/st0 /data
```

Additional resources

- **xfsdump(8)** man page.

18.3. ADDITIONAL RESOURCES

- **xfsdump(8)** man page

CHAPTER 19. RESTORING AN XFS FILE SYSTEM FROM BACKUP

As a system administrator, you can use the **xfsrestore** utility to restore XFS backup created with the **xfsdump** utility and stored in a file or on a tape.

19.1. FEATURES OF RESTORING XFS FROM BACKUP

The **xfsrestore** utility restores file systems from backups produced by **xfsdump**. The **xfsrestore** utility has two modes:

- The **simple** mode enables users to restore an entire file system from a level 0 dump. This is the default mode.
- The **cumulative** mode enables file system restoration from an incremental backup: that is, level 1 to level 9.

A unique *session ID* or *session label* identifies each backup. Restoring a backup from a tape containing multiple backups requires its corresponding session ID or label.

To extract, add, or delete specific files from a backup, enter the **xfsrestore** interactive mode. The interactive mode provides a set of commands to manipulate the backup files.

Additional resources

- **xfsrestore(8)** man page.

19.2. RESTORING AN XFS FILE SYSTEM FROM BACKUP WITH XFSRESTORE

This procedure describes how to restore the content of an XFS file system from a file or tape backup.

Prerequisites

- A file or tape backup of XFS file systems, as described in [Backing up an XFS file system](#).
- A storage device where you can restore the backup.

Procedure

- The command to restore the backup varies depending on whether you are restoring from a full backup or an incremental one, or are restoring multiple backups from a single tape device:

```
# xfsrestore [-r] [-S session-id] [-L session-label] [-i]  
-f backup-location restoration-path
```

- Replace *backup-location* with the location of the backup. This can be a regular file, a tape drive, or a remote tape device. For example, **/backup-files/Data.xfsdump** for a file or **/dev/st0** for a tape drive.
- Replace *restoration-path* with the path to the directory where you want to restore the file system. For example, **/mnt/data/**.

- To restore a file system from an incremental (level 1 to level 9) backup, add the **-r** option.
- To restore a backup from a tape device that contains multiple backups, specify the backup using the **-S** or **-L** options.
The **-S** option lets you choose a backup by its session ID, while the **-L** option lets you choose by the session label. To obtain the session ID and session labels, use the **xfsrestore -l** command.

Replace *session-id* with the session ID of the backup. For example, **b74a3586-e52e-4a4a-8775-c3334fa8ea2c**. Replace *session-label* with the session label of the backup. For example, **my_backup_session_label**.

- To use **xfsrestore** interactively, use the **-i** option.
The interactive dialog begins after **xfsrestore** finishes reading the specified device. Available commands in the interactive **xfsrestore** shell include **cd**, **ls**, **add**, **delete**, and **extract**; for a complete list of commands, use the **help** command.

Example 19.1. Restoring Multiple XFS File Systems

- To restore the XFS backup files and save their content into directories under **/mnt/**:

```
# xfsrestore -f /backup-files/boot.xfsdump /mnt/boot/
# xfsrestore -f /backup-files/data.xfsdump /mnt/data/
```

- To restore from a tape device containing multiple backups, specify each backup by its session label or session ID:

```
# xfsrestore -L "backup_boot" -f /dev/st0 /mnt/boot/
# xfsrestore -S "45e9af35-efd2-4244-87bc-4762e476cbab" \
-f /dev/st0 /mnt/data/
```

Additional resources

- **xfsrestore(8)** man page.

19.3. INFORMATIONAL MESSAGES WHEN RESTORING AN XFS BACKUP FROM A TAPE

When restoring a backup from a tape with backups from multiple file systems, the **xfsrestore** utility might issue messages. The messages inform you whether a match of the requested backup has been found when **xfsrestore** examines each backup on the tape in sequential order. For example:

```
xfsrestore: preparing drive
xfsrestore: examining media file 0
xfsrestore: inventory session uuid (8590224e-3c93-469c-a311-fc8f23029b2a) does not match the
media header's session uuid (7eda9f86-f1e9-4dfd-b1d4-c50467912408)
xfsrestore: examining media file 1
xfsrestore: inventory session uuid (8590224e-3c93-469c-a311-fc8f23029b2a) does not match the
media header's session uuid (7eda9f86-f1e9-4dfd-b1d4-c50467912408)
[...]
```

The informational messages keep appearing until the matching backup is found.

19.4. ADDITIONAL RESOURCES

- **xfsrestore(8)** man page

CHAPTER 20. INCREASING THE SIZE OF AN XFS FILE SYSTEM

As a system administrator, you can increase the size of an XFS file system to make a complete use of a larger storage capacity.



IMPORTANT

It is not currently possible to decrease the size of XFS file systems.

20.1. INCREASING THE SIZE OF AN XFS FILE SYSTEM WITH XFS_GROWFS

This procedure describes how to grow an XFS file system using the **xfs_growfs** utility.

Prerequisites

- Ensure that the underlying block device is of an appropriate size to hold the resized file system later. Use the appropriate resizing methods for the affected block device.
- Mount the XFS file system.

Procedure

- While the XFS file system is mounted, use the **xfs_growfs** utility to increase its size:

```
# xfs_growfs file-system -D new-size
```

- Replace *file-system* with the mount point of the XFS file system.
- With the **-D** option, replace *new-size* with the desired new size of the file system specified in the number of file system blocks.

To find out the block size in kB of a given XFS file system, use the **xfs_info** utility:

```
# xfs_info block-device
...
data    =          bsize=4096
...
```

- Without the **-D** option, **xfs_growfs** grows the file system to the maximum size supported by the underlying device.

Additional resources

- **xfs_growfs(8)** man page.

CHAPTER 21. CONFIGURING XFS ERROR BEHAVIOR

You can configure how an XFS file system behaves when it encounters different I/O errors.

21.1. CONFIGURABLE ERROR HANDLING IN XFS

The XFS file system responds in one of the following ways when an error occurs during an I/O operation:

- XFS repeatedly retries the I/O operation until the operation succeeds or XFS reaches a set limit. The limit is based either on a maximum number of retries or a maximum time for retries.
- XFS considers the error permanent and stops the operation on the file system.

You can configure how XFS reacts to the following error conditions:

EIO

Error when reading or writing

ENOSPC

No space left on the device

ENODEV

Device cannot be found

You can set the maximum number of retries and the maximum time in seconds until XFS considers an error permanent. XFS stops retrying the operation when it reaches either of the limits.

You can also configure XFS so that when unmounting a file system, XFS immediately cancels the retries regardless of any other configuration. This configuration enables the unmount operation to succeed despite persistent errors.

Default behavior

The default behavior for each XFS error condition depends on the error context. Some XFS errors such as **ENODEV** are considered to be fatal and unrecoverable, regardless of the retry count. Their default retry limit is 0.

21.2. CONFIGURATION FILES FOR SPECIFIC AND UNDEFINED XFS ERROR CONDITIONS

The following directories store configuration files that control XFS error behavior for different error conditions:

/sys/fs/xfs/device/error/metadata/EIO/

For the **EIO** error condition

/sys/fs/xfs/device/error/metadata/ENODEV/

For the **ENODEV** error condition

/sys/fs/xfs/device/error/metadata/ENOSPC/

For the **ENOSPC** error condition

/sys/fs/xfs/device/error/default/

Common configuration for all other, undefined error conditions

Each directory contains the following configuration files for configuring retry limits:

max_retries

Controls the maximum number of times that XFS retries the operation.

retry_timeout_seconds

Specifies the time limit in seconds after which XFS stops retrying the operation.

21.3. SETTING XFS BEHAVIOR FOR SPECIFIC CONDITIONS

This procedure configures how XFS reacts to specific error conditions.

Procedure

- Set the maximum number of retries, the retry time limit, or both:
 - To set the maximum number of retries, write the desired number to the **max_retries** file:

```
# echo value > /sys/fs/xfs/device/error/metadata/condition/max_retries
```

- To set the time limit, write the desired number of seconds to the **retry_timeout_seconds** file:

```
# echo value > /sys/fs/xfs/device/error/metadata/condition/retry_timeout_second
```

value is a number between -1 and the maximum possible value of the C signed integer type. This is 2147483647 on 64-bit Linux.

In both limits, the value **-1** is used for continuous retries and **0** to stop immediately.

device is the name of the device, as found in the **/dev/** directory; for example, **sda**.

21.4. SETTING XFS BEHAVIOR FOR UNDEFINED CONDITIONS

This procedure configures how XFS reacts to all undefined error conditions, which share a common configuration.

Procedure

- Set the maximum number of retries, the retry time limit, or both:
 - To set the maximum number of retries, write the desired number to the **max_retries** file:

```
# echo value > /sys/fs/xfs/device/error/metadata/default/max_retries
```

- To set the time limit, write the desired number of seconds to the **retry_timeout_seconds** file:

```
# echo value > /sys/fs/xfs/device/error/metadata/default/retry_timeout_seconds
```

value is a number between -1 and the maximum possible value of the C signed integer type. This is 2147483647 on 64-bit Linux.

In both limits, the value **-1** is used for continuous retries and **0** to stop immediately.

device is the name of the device, as found in the **/dev/** directory; for example, **sda**.

21.5. SETTING THE XFS UNMOUNT BEHAVIOR

This procedure configures how XFS reacts to error conditions when unmounting the file system.

If you set the **fail_at_unmount** option in the file system, it overrides all other error configurations during unmount, and immediately unmounts the file system without retrying the I/O operation. This allows the unmount operation to succeed even in case of persistent errors.



WARNING

You cannot change the **fail_at_unmount** value after the unmount process starts, because the unmount process removes the configuration files from the **sysfs** interface for the respective file system. You must configure the unmount behavior before the file system starts unmounting.

Procedure

- Enable or disable the **fail_at_unmount** option:
 - To cancel retrying all operations when the file system unmounts, enable the option:

```
# echo 1 > /sys/fs/xfs/device/error/fail_at_unmount
```

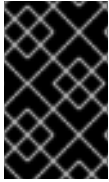
- To respect the **max_retries** and **retry_timeout_seconds** retry limits when the file system unmounts, disable the option:

```
# echo 0 > /sys/fs/xfs/device/error/fail_at_unmount
```

device is the name of the device, as found in the **/dev/** directory; for example, **sda**.

CHAPTER 22. CHECKING AND REPAIRING A FILE SYSTEM

RHEL provides file system administration utilities which are capable of checking and repairing file systems. These tools are often referred to as **fsck** tools, where **fsck** is a shortened version of *file system check*. In most cases, these utilities are run automatically during system boot, if needed, but can also be manually invoked if required.



IMPORTANT

File system checkers guarantee only metadata consistency across the file system. They have no awareness of the actual data contained within the file system and are not data recovery tools.

22.1. SCENARIOS THAT REQUIRE A FILE SYSTEM CHECK

The relevant **fsck** tools can be used to check your system if any of the following occurs:

- System fails to boot
- Files on a specific disk become corrupt
- The file system shuts down or changes to read-only due to inconsistencies
- A file on the file system is inaccessible

File system inconsistencies can occur for various reasons, including but not limited to hardware errors, storage administration errors, and software bugs.



IMPORTANT

File system check tools cannot repair hardware problems. A file system must be fully readable and writable if repair is to operate successfully. If a file system was corrupted due to a hardware error, the file system must first be moved to a good disk, for example with the **dd(8)** utility.

For journaling file systems, all that is normally required at boot time is to replay the journal if required and this is usually a very short operation.

However, if a file system inconsistency or corruption occurs, even for journaling file systems, then the file system checker must be used to repair the file system.



IMPORTANT

It is possible to disable file system check at boot by setting the sixth field in **/etc/fstab** to **0**. However, Red Hat does not recommend doing so unless you are having issues with **fsck** at boot time, for example with extremely large or remote file systems.

Additional resources

- **fstab(5)** man page.
- **fsck(8)** man page.
- **dd(8)** man page.

22.2. POTENTIAL SIDE EFFECTS OF RUNNING FSCK

Generally, running the file system check and repair tool can be expected to automatically repair at least some of the inconsistencies it finds. In some cases, the following issues can arise:

- Severely damaged inodes or directories may be discarded if they cannot be repaired.
- Significant changes to the file system may occur.

To ensure that unexpected or undesirable changes are not permanently made, ensure you follow any precautionary steps outlined in the procedure.

22.3. ERROR-HANDLING MECHANISMS IN XFS

This section describes how XFS handles various kinds of errors in the file system.

Unclean unmounts

Journaling maintains a transactional record of metadata changes that happen on the file system.

In the event of a system crash, power failure, or other unclean unmount, XFS uses the journal (also called log) to recover the file system. The kernel performs journal recovery when mounting the XFS file system.

Corruption

In this context, *corruption* means errors on the file system caused by, for example:

- Hardware faults
- Bugs in storage firmware, device drivers, the software stack, or the file system itself
- Problems that cause parts of the file system to be overwritten by something outside of the file system

When XFS detects corruption in the file system or the file-system metadata, it may shut down the file system and report the incident in the system log. Note that if the corruption occurred on the file system hosting the **/var** directory, these logs will not be available after a reboot.

Example 22.1. System log entry reporting an XFS corruption

```
# dmesg --notime | tail -15
```

```
XFS (loop0): Mounting V5 Filesystem
```

```
XFS (loop0): Metadata CRC error detected at xfs_agi_read_verify+0xcb/0xf0 [xfs], xfs_agi block 0x2
```

```
XFS (loop0): Unmount and run xfs_repair
```

```
XFS (loop0): First 128 bytes of corrupted metadata buffer:
```

```
00000000027b3b56: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
```

```
0000000005f9abc7a: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
```

```
0000000005b0aef35: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
```

```
000000000da9d2ded: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
```

```
0000000001e265b07: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
```

```
0000000006a40df69: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
```

```
000000000b272907: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
```

```
000000000e484aac5: 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 00 .....
```

```
XFS (loop0): metadata I/O error in "xfs_trans_read_buf_map" at daddr 0x2 len 1 error 74
XFS (loop0): xfs_imap_lookup: xfs_ialloc_read_agi() returned error -117, agno 0
XFS (loop0): Failed to read root inode 0x80, error 11
```

User-space utilities usually report the *Input/output error* message when trying to access a corrupted XFS file system. Mounting an XFS file system with a corrupted log results in a failed mount and the following error message:

```
mount: /mount-point: mount(2) system call failed: Structure needs cleaning.
```

You must manually use the **xfs_repair** utility to repair the corruption.

Additional resources

- **xfs_repair(8)** man page.

22.4. CHECKING AN XFS FILE SYSTEM WITH XFS_REPAIR

This procedure performs a read-only check of an XFS file system using the **xfs_repair** utility. You must manually use the **xfs_repair** utility to repair any corruption. Unlike other file system repair utilities, **xfs_repair** does not run at boot time, even when an XFS file system was not cleanly unmounted. In the event of an unclean unmount, XFS simply replays the log at mount time, ensuring a consistent file system; **xfs_repair** cannot repair an XFS file system with a dirty log without remounting it first.



NOTE

Although an **fsck.xfs** binary is present in the **xfsprogs** package, this is present only to satisfy **initscripts** that look for an **fsck.file** system binary at boot time. **fsck.xfs** immediately exits with an exit code of 0.

Procedure

1. Replay the log by mounting and unmounting the file system:

```
# mount file-system
# umount file-system
```



NOTE

If the mount fails with a structure needs cleaning error, the log is corrupted and cannot be replayed. The dry run should discover and report more on-disk corruption as a result.

2. Use the **xfs_repair** utility to perform a dry run to check the file system. Any errors are printed and an indication of the actions that would be taken, without modifying the file system.

```
# xfs_repair -n block-device
```

3. Mount the file system:

```
# mount file-system
```

Additional resources

- **xfs_repair(8)** man page.
- **xfs_metadump(8)** man page.

22.5. REPAIRING AN XFS FILE SYSTEM WITH XFS_REPAIR

This procedure repairs a corrupted XFS file system using the **xfs_repair** utility.

Procedure

1. Create a metadata image prior to repair for diagnostic or testing purposes using the **xfs_metadump** utility. A pre-repair file system metadata image can be useful for support investigations if the corruption is due to a software bug. Patterns of corruption present in the pre-repair image can aid in root-cause analysis.
 - Use the **xfs_metadump** debugging tool to copy the metadata from an XFS file system to a file. The resulting **metadump** file can be compressed using standard compression utilities to reduce the file size if large **metadump** files need to be sent to support.

```
# xfs_metadump block-device metadump-file
```

2. Replay the log by remounting the file system:

```
# mount file-system  
# umount file-system
```

3. Use the **xfs_repair** utility to repair the unmounted file system:

- If the mount succeeded, no additional options are required:

```
# xfs_repair block-device
```

- If the mount failed with the *Structure needs cleaning* error, the log is corrupted and cannot be replayed. Use the **-L** option (*force log zeroing*) to clear the log:



WARNING

This command causes all metadata updates in progress at the time of the crash to be lost, which might cause significant file system damage and data loss. This should be used only as a last resort if the log cannot be replayed.

```
# xfs_repair -L block-device
```

4. Mount the file system:

```
# mount file-system
```


Additional resources

- **xfs_repair(8)** man page.

22.6. ERROR HANDLING MECHANISMS IN EXT2, EXT3, AND EXT4

The ext2, ext3, and ext4 file systems use the **e2fsck** utility to perform file system checks and repairs. The file names **fsck.ext2**, **fsck.ext3**, and **fsck.ext4** are hardlinks to the **e2fsck** utility. These binaries are run automatically at boot time and their behavior differs based on the file system being checked and the state of the file system.

A full file system check and repair is invoked for ext2, which is not a metadata journaling file system, and for ext4 file systems without a journal.

For ext3 and ext4 file systems with metadata journaling, the journal is replayed in userspace and the utility exits. This is the default action because journal replay ensures a consistent file system after a crash.

If these file systems encounter metadata inconsistencies while mounted, they record this fact in the file system superblock. If **e2fsck** finds that a file system is marked with such an error, **e2fsck** performs a full check after replaying the journal (if present).

Additional resources

- **fsck(8)** man page.
- **e2fsck(8)** man page.

22.7. CHECKING AN EXT2, EXT3, OR EXT4 FILE SYSTEM WITH E2FSCK

This procedure checks an ext2, ext3, or ext4 file system using the **e2fsck** utility.

Procedure

1. Replay the log by remounting the file system:

```
# mount file-system
# umount file-system
```

2. Perform a dry run to check the file system.

```
# e2fsck -n block-device
```



NOTE

Any errors are printed and an indication of the actions that would be taken, without modifying the file system. Later phases of consistency checking may print extra errors as it discovers inconsistencies which would have been fixed in early phases if it were running in repair mode.

Additional resources

- **e2image(8)** man page.
- **e2fsck(8)** man page.

22.8. REPAIRING AN EXT2, EXT3, OR EXT4 FILE SYSTEM WITH E2FSCK

This procedure repairs a corrupted ext2, ext3, or ext4 file system using the **e2fsck** utility.

Procedure

1. Save a file system image for support investigations. A pre-repair file system metadata image can be useful for support investigations if the corruption is due to a software bug. Patterns of corruption present in the pre-repair image can aid in root-cause analysis.



NOTE

Severely damaged file systems may cause problems with metadata image creation.

- If you are creating the image for testing purposes, use the **-r** option to create a sparse file of the same size as the file system itself. **e2fsck** can then operate directly on the resulting file.

```
# e2image -r block-device image-file
```

- If you are creating the image to be archived or provided for diagnostic, use the **-Q** option, which creates a more compact file format suitable for transfer.

```
# e2image -Q block-device image-file
```

2. Replay the log by remounting the file system:

```
# mount file-system  
# umount file-system
```

3. Automatically repair the file system. If user intervention is required, **e2fsck** indicates the unfixed problem in its output and reflects this status in the exit code.

```
# e2fsck -p block-device
```

Additional resources

- **e2image(8)** man page.
- **e2fsck(8)** man page.

CHAPTER 23. MOUNTING FILE SYSTEMS

As a system administrator, you can mount file systems on your system to access data on them.

23.1. THE LINUX MOUNT MECHANISM

This section explains basic concepts of mounting file systems on Linux.

On Linux, UNIX, and similar operating systems, file systems on different partitions and removable devices (CDs, DVDs, or USB flash drives for example) can be attached to a certain point (the mount point) in the directory tree, and then detached again. While a file system is mounted on a directory, the original content of the directory is not accessible.

Note that Linux does not prevent you from mounting a file system to a directory with a file system already attached to it.

When mounting, you can identify the device by:

- a universally unique identifier (UUID): for example, **UUID=34795a28-ca6d-4fd8-a347-73671d0c19cb**
- a volume label: for example, **LABEL=home**
- a full path to a non-persistent block device: for example, **/dev/sda3**

When you mount a file system using the **mount** command without all required information, that is without the device name, the target directory, or the file system type, the **mount** utility reads the content of the **/etc/fstab** file to check if the given file system is listed there. The **/etc/fstab** file contains a list of device names and the directories in which the selected file systems are set to be mounted as well as the file system type and mount options. Therefore, when mounting a file system that is specified in **/etc/fstab**, the following command syntax is sufficient:

- Mounting by the mount point:

```
# mount directory
```

- Mounting by the block device:

```
# mount device
```

Additional resources

- **mount(8)** man page
- [How to list persistent naming attributes such as the UUID](#) .

23.2. LISTING CURRENTLY MOUNTED FILE SYSTEMS

This procedure describes how to list all currently mounted file systems on the command line.

Procedure

- To list all mounted file systems, use the **findmnt** utility:

```
$ findmnt
```

- To limit the listed file systems only to a certain file system type, add the **--types** option:

```
$ findmnt --types fs-type
```

For example:

Example 23.1. Listing only XFS file systems

```
$ findmnt --types xfs
```

TARGET	SOURCE	FSTYPE	OPTIONS
/	/dev/mapper/luks-5564ed00-6aac-4406-bfb4-c59bf5de48b5	xfs	rw,relatime
└-/boot	/dev/sda1	xfs	rw,relatime
└-/home	/dev/mapper/luks-9d185660-7537-414d-b727-d92ea036051e	xfs	rw,relatime

Additional resources

- **findmnt(8)** man page

23.3. MOUNTING A FILE SYSTEM WITH MOUNT

This procedure describes how to mount a file system using the **mount** utility.

Prerequisites

- Make sure that no file system is already mounted on your chosen mount point:

```
$ findmnt mount-point
```

Procedure

1. To attach a certain file system, use the **mount** utility:

```
# mount device mount-point
```

Example 23.2. Mounting an XFS file system

For example, to mount a local XFS file system identified by UUID:

```
# mount UUID=ea74bbec-536d-490c-b8d9-5b40bbd7545b /mnt/data
```

2. If **mount** cannot recognize the file system type automatically, specify it using the **--types** option:

```
# mount --types type device mount-point
```

Example 23.3. Mounting an NFS file system

For example, to mount a remote NFS file system:

```
# mount --types nfs4 host:/remote-export /mnt/nfs
```

Additional resources

- **mount(8)** man page

23.4. MOVING A MOUNT POINT

This procedure describes how to change the mount point of a mounted file system to a different directory.

Procedure

1. To change the directory in which a file system is mounted:

```
# mount --move old-directory new-directory
```

Example 23.4. Moving a home file system

For example, to move the file system mounted in the **/mnt/userdirs/** directory to the **/home/** mount point:

```
# mount --move /mnt/userdirs /home
```

2. Verify that the file system has been moved as expected:

```
$ findmnt
$ ls old-directory
$ ls new-directory
```

Additional resources

- **mount(8)** man page

23.5. UNMOUNTING A FILE SYSTEM WITH Umount

This procedure describes how to unmount a file system using the **umount** utility.

Procedure

1. Try unmounting the file system using either of the following commands:

- By mount point:

```
# umount mount-point
```

- By device:

```
# umount device
```

If the command fails with an error similar to the following, it means that the file system is in use because of a process is using resources on it:

```
umount: /run/media/user/FlashDrive: target is busy.
```

- If the file system is in use, use the **fuser** utility to determine which processes are accessing it. For example:

```
$ fuser --mount /run/media/user/FlashDrive

/run/media/user/FlashDrive: 18351
```

Afterwards, terminate the processes using the file system and try unmounting it again.

23.6. COMMON MOUNT OPTIONS

The following table lists the most common options of the **mount** utility. You can apply these mount options using the following syntax:

```
# mount --options option1,option2,option3 device mount-point
```

Table 23.1. Common mount options

Option	Description
async	Enables asynchronous input and output operations on the file system.
auto	Enables the file system to be mounted automatically using the mount -a command.
defaults	Provides an alias for the async,auto,dev,exec,nouser,rw,suid options.
exec	Allows the execution of binary files on the particular file system.
loop	Mounts an image as a loop device.
noauto	Default behavior disables the automatic mount of the file system using the mount -a command.
noexec	Disallows the execution of binary files on the particular file system.
nouser	Disallows an ordinary user (that is, other than root) to mount and unmount the file system.
remount	Remounts the file system in case it is already mounted.
ro	Mounts the file system for reading only.

Option	Description
rw	Mounts the file system for both reading and writing.
user	Allows an ordinary user (that is, other than root) to mount and unmount the file system.

CHAPTER 24. SHARING A MOUNT ON MULTIPLE MOUNT POINTS

As a system administrator, you can duplicate mount points to make the file systems accessible from multiple directories.

24.1. TYPES OF SHARED MOUNTS

There are multiple types of shared mounts that you can use. The difference between them is what happens when you mount another file system under one of the shared mount points. The shared mounts are implemented using the *shared subtrees* functionality.

The following mount types are available:

private

This type does not receive or forward any propagation events.

When you mount another file system under either the duplicate or the original mount point, it is not reflected in the other.

shared

This type creates an exact replica of a given mount point.

When a mount point is marked as a **shared** mount, any mount within the original mount point is reflected in it, and vice versa.

This is the default mount type of the root file system.

slave

This type creates a limited duplicate of a given mount point.

When a mount point is marked as a **slave** mount, any mount within the original mount point is reflected in it, but no mount within a **slave** mount is reflected in its original.

unbindable

This type prevents the given mount point from being duplicated whatsoever.

Additional resources

- [The *Shared subtrees* article on Linux Weekly News.](#)

24.2. CREATING A PRIVATE MOUNT POINT DUPLICATE

This procedure duplicates a mount point as a private mount. File systems that you later mount under the duplicate or the original mount point are not reflected in the other.

Procedure

1. Create a virtual file system (VFS) node from the original mount point:

```
# mount --bind original-dir original-dir
```

2. Mark the original mount point as private:


```
# mount --make-private original-dir
```

Alternatively, to change the mount type for the selected mount point and all mount points under it, use the **--make-rprivate** option instead of **--make-private**.

3. Create the duplicate:

```
# mount --bind original-dir duplicate-dir
```

Example 24.1. Duplicating `/media` into `/mnt` as a private mount point

1. Create a VFS node from the `/media` directory:

```
# mount --bind /media /media
```

2. Mark the `/media` directory as private:

```
# mount --make-private /media
```

3. Create its duplicate in `/mnt`:

```
# mount --bind /media /mnt
```

4. It is now possible to verify that `/media` and `/mnt` share content but none of the mounts within `/media` appear in `/mnt`. For example, if the CD-ROM drive contains non-empty media and the `/media/cdrom/` directory exists, use:

```
# mount /dev/cdrom /media/cdrom
# ls /media/cdrom
EFI  GPL  isolinux  LiveOS
# ls /mnt/cdrom
#
```

5. It is also possible to verify that file systems mounted in the `/mnt` directory are not reflected in `/media`. For instance, if a non-empty USB flash drive that uses the `/dev/sdc1` device is plugged in and the `/mnt/flashdisk/` directory is present, use:

```
# mount /dev/sdc1 /mnt/flashdisk
# ls /media/flashdisk
# ls /mnt/flashdisk
en-US  publican.cfg
```

Additional resources

- **mount(8)** man page

24.3. CREATING A SHARED MOUNT POINT DUPLICATE

This procedure duplicates a mount point as a shared mount. File systems that you later mount under the original directory or the duplicate are always reflected in the other.

Procedure

1. Create a virtual file system (VFS) node from the original mount point:

```
# mount --bind original-dir original-dir
```

2. Mark the original mount point as shared:

```
# mount --make-shared original-dir
```

Alternatively, to change the mount type for the selected mount point and all mount points under it, use the **--make-rshared** option instead of **--make-shared**.

3. Create the duplicate:

```
# mount --bind original-dir duplicate-dir
```

Example 24.2. Duplicating /media into /mnt as a shared mount point

To make the **/media** and **/mnt** directories share the same content:

1. Create a VFS node from the **/media** directory:

```
# mount --bind /media /media
```

2. Mark the **/media** directory as shared:

```
# mount --make-shared /media
```

3. Create its duplicate in **/mnt**:

```
# mount --bind /media /mnt
```

4. It is now possible to verify that a mount within **/media** also appears in **/mnt**. For example, if the CD-ROM drive contains non-empty media and the **/media/cdrom/** directory exists, use:

```
# mount /dev/cdrom /media/cdrom
# ls /media/cdrom
EFI  GPL  isolinux  LiveOS
# ls /mnt/cdrom
EFI  GPL  isolinux  LiveOS
```

5. Similarly, it is possible to verify that any file system mounted in the **/mnt** directory is reflected in **/media**. For instance, if a non-empty USB flash drive that uses the **/dev/sdc1** device is plugged in and the **/mnt/flashdisk/** directory is present, use:

```
# mount /dev/sdc1 /mnt/flashdisk
# ls /media/flashdisk
en-US  publican.cfg
# ls /mnt/flashdisk
en-US  publican.cfg
```

Additional resources

- **mount(8)** man page

24.4. CREATING A SLAVE MOUNT POINT DUPLICATE

This procedure duplicates a mount point as a **slave** mount type. File systems that you later mount under the original mount point are reflected in the duplicate but not the other way around.

Procedure

1. Create a virtual file system (VFS) node from the original mount point:

```
# mount --bind original-dir original-dir
```

2. Mark the original mount point as shared:

```
# mount --make-shared original-dir
```

Alternatively, to change the mount type for the selected mount point and all mount points under it, use the **--make-rshared** option instead of **--make-shared**.

3. Create the duplicate and mark it as the **slave** type:

```
# mount --bind original-dir duplicate-dir
# mount --make-slave duplicate-dir
```

Example 24.3. Duplicating **/media** into **/mnt** as a slave mount point

This example shows how to get the content of the **/media** directory to appear in **/mnt** as well, but without any mounts in the **/mnt** directory to be reflected in **/media**.

1. Create a VFS node from the **/media** directory:

```
# mount --bind /media /media
```

2. Mark the **/media** directory as shared:

```
# mount --make-shared /media
```

3. Create its duplicate in **/mnt** and mark it as **slave**:

```
# mount --bind /media /mnt
# mount --make-slave /mnt
```

4. Verify that a mount within **/media** also appears in **/mnt**. For example, if the CD-ROM drive contains non-empty media and the **/media/cdrom/** directory exists, use:

```
# mount /dev/cdrom /media/cdrom
# ls /media/cdrom
EFI GPL isolinux LiveOS
# ls /mnt/cdrom
EFI GPL isolinux LiveOS
```

5. Also verify that file systems mounted in the **/mnt** directory are not reflected in **/media**. For instance, if a non-empty USB flash drive that uses the **/dev/sdc1** device is plugged in and the **/mnt/flashdisk/** directory is present, use:

```
# mount /dev/sdc1 /mnt/flashdisk
# ls /media/flashdisk
# ls /mnt/flashdisk
en-US publican.cfg
```

Additional resources

- **mount(8)** man page

24.5. PREVENTING A MOUNT POINT FROM BEING DUPLICATED

This procedure marks a mount point as unbindable so that it is not possible to duplicate it in another mount point.

Procedure

- To change the type of a mount point to an unbindable mount, use:

```
# mount --bind mount-point mount-point
# mount --make-unbindable mount-point
```

Alternatively, to change the mount type for the selected mount point and all mount points under it, use the **--make-runbindable** option instead of **--make-unbindable**.

Any subsequent attempt to make a duplicate of this mount fails with the following error:

```
# mount --bind mount-point duplicate-dir

mount: wrong fs type, bad option, bad superblock on mount-point,
missing codepage or helper program, or other error
In some cases useful info is found in syslog - try
dmesg | tail or so
```

Example 24.4. Preventing **/media** from being duplicated

- To prevent the **/media** directory from being shared, use:

```
# mount --bind /media /media
# mount --make-unbindable /media
```

Additional resources

- **mount(8)** man page

CHAPTER 25. PERSISTENTLY MOUNTING FILE SYSTEMS

As a system administrator, you can persistently mount file systems to configure non-removable storage.

25.1. THE /ETC/FSTAB FILE

Use the **/etc/fstab** configuration file to control persistent mount points of file systems. Each line in the **/etc/fstab** file defines a mount point of a file system.

It includes six fields separated by white space:

1. The block device identified by a persistent attribute or a path in the **/dev** directory.
2. The directory where the device will be mounted.
3. The file system on the device.
4. Mount options for the file system, which includes the **defaults** option to mount the partition at boot time with default options. The mount option field also recognizes the **systemd** mount unit options in the **x-systemd.option** format.
5. Backup option for the **dump** utility.
6. Check order for the **fsck** utility.



NOTE

The **systemd-fstab-generator** dynamically converts the entries from the **/etc/fstab** file to the **systemd-mount** units. The **systemd** auto mounts LVM volumes from **/etc/fstab** during manual activation unless the **systemd-mount** unit is masked.



NOTE

The **dump** utility used for backup of file systems has been removed in RHEL 9, and is available in the EPEL 9 repository.

Example 25.1. The **/boot** file system in **/etc/fstab**

Block device	Mount point	File system	Options	Backup	Check
UUID=ea74bbec-536d-490c-b8d9-5b40bbd7545b	/boot	xfs	defaults	0	0

The **systemd** service automatically generates mount units from entries in **/etc/fstab**.

Additional resources

- **fstab(5)** and **systemd.mount(5)** man pages

25.2. ADDING A FILE SYSTEM TO /ETC/FSTAB

This procedure describes how to configure persistent mount point for a file system in the **/etc/fstab** configuration file.

Procedure

1. Find out the UUID attribute of the file system:

```
$ lsblk --fs storage-device
```

For example:

Example 25.2. Viewing the UUID of a partition

```
$ lsblk --fs /dev/sda1
```

NAME	FSTYPE	LABEL	UUID	MOUNTPOINT
sda1	xf	Boot	ea74bbec-536d-490c-b8d9-5b40bbd7545b	/boot

2. If the mount point directory does not exist, create it:

```
# mkdir --parents mount-point
```

3. As root, edit the **/etc/fstab** file and add a line for the file system, identified by the UUID.
For example:

Example 25.3. The /boot mount point in /etc/fstab

```
UUID=ea74bbec-536d-490c-b8d9-5b40bbd7545b /boot xfs defaults 0 0
```

4. Regenerate mount units so that your system registers the new configuration:

```
# systemctl daemon-reload
```

5. Try mounting the file system to verify that the configuration works:

```
# mount mount-point
```

Additional resources

- [Overview of persistent naming attributes.](#)

CHAPTER 26. MOUNTING FILE SYSTEMS ON DEMAND

As a system administrator, you can configure file systems, such as NFS, to mount automatically on demand.

26.1. THE AUTOFS SERVICE

This section explains the benefits and basic concepts of the **autofs** service, used to mount file systems on demand.

One drawback of permanent mounting using the **/etc/fstab** configuration is that, regardless of how infrequently a user accesses the mounted file system, the system must dedicate resources to keep the mounted file system in place. This might affect system performance when, for example, the system is maintaining NFS mounts to many systems at one time.

An alternative to **/etc/fstab** is to use the kernel-based **autofs** service. It consists of the following components:

- A kernel module that implements a file system, and
- A user-space service that performs all of the other functions.

The **autofs** service can mount and unmount file systems automatically (on-demand), therefore saving system resources. It can be used to mount file systems such as NFS, AFS, SMBFS, CIFS, and local file systems.

Additional resources

- The **autofs(8)** man page.

26.2. THE AUTOFS CONFIGURATION FILES

This section describes the usage and syntax of configuration files used by the **autofs** service.

The master map file

The **autofs** service uses **/etc/auto.master** (master map) as its default primary configuration file. This can be changed to use another supported network source and name using the **autofs** configuration in the **/etc/autofs.conf** configuration file in conjunction with the Name Service Switch (NSS) mechanism.

All on-demand mount points must be configured in the master map. Mount point, host name, exported directory, and options can all be specified in a set of files (or other supported network sources) rather than configuring them manually for each host.

The master map file lists mount points controlled by **autofs**, and their corresponding configuration files or network sources known as automount maps. The format of the master map is as follows:

```
mount-point map-name options
```

The variables used in this format are:

mount-point

The **autofs** mount point; for example, **/mnt/data**.

map-file

The map source file, which contains a list of mount points and the file system location from which those mount points should be mounted.

options

If supplied, these apply to all entries in the given map, if they do not themselves have options specified.

Example 26.1. The `/etc/auto.master` file

The following is a sample line from **`/etc/auto.master`** file:

```
/mnt/data /etc/auto.data
```

Map files

Map files configure the properties of individual on-demand mount points.

The automounter creates the directories if they do not exist. If the directories exist before the automounter was started, the automounter will not remove them when it exits. If a timeout is specified, the directory is automatically unmounted if the directory is not accessed for the timeout period.

The general format of maps is similar to the master map. However, the options field appears between the mount point and the location instead of at the end of the entry as in the master map:

```
mount-point options location
```

The variables used in this format are:

mount-point

This refers to the **autofs** mount point. This can be a single directory name for an indirect mount or the full path of the mount point for direct mounts. Each direct and indirect map entry key (*mount-point*) can be followed by a space separated list of offset directories (subdirectory names each beginning with */*) making them what is known as a multi-mount entry.

options

When supplied, these options are appended to the master map entry options, if any, or used instead of the master map options if the configuration entry **append_options** is set to **no**.

location

This refers to the file system location such as a local file system path (preceded with the Sun map format escape character **:** for map names beginning with */*), an NFS file system or other valid file system location.

Example 26.2. A map file

The following is a sample from a map file; for example, **`/etc/auto.misc`**:

```
payroll -fstype=nfs4 personnel:/exports/payroll
sales -fstype=xfss /dev/hda4
```

The first column in the map file indicates the **autofs** mount point: **sales** and **payroll** from the server called **personnel**. The second column indicates the options for the **autofs** mount. The third column indicates the source of the mount.

Following the given configuration, the **autofs** mount points will be **/home/payroll** and **/home/sales**. The **-fstype=** option is often omitted and is not needed if the file system is NFS, including mounts for NFSv4 if the system default is NFSv4 for NFS mounts.

Using the given configuration, if a process requires access to an **autofs** unmounted directory such as **/home/payroll/2006/July.sxc**, the **autofs** service automatically mounts the directory.

The amd map format

The **autofs** service recognizes map configuration in the **amd** format as well. This is useful if you want to reuse existing automounter configuration written for the **am-utils** service, which has been removed from Red Hat Enterprise Linux.

However, Red Hat recommends using the simpler **autofs** format described in the previous sections.

Additional resources

- **autofs(5)** man page
- **autofs.conf(5)** man page
- **auto.master(5)** man page
- **/usr/share/doc/autofs/README.amd-maps** file

26.3. CONFIGURING AUTOFS MOUNT POINTS

This procedure describes how to configure on-demand mount points using the **autofs** service.

Prerequisites

- Install the **autofs** package:

```
# dnf install autofs
```

- Start and enable the **autofs** service:

```
# systemctl enable --now autofs
```

Procedure

1. Create a map file for the on-demand mount point, located at **/etc/auto.*identifier***. Replace *identifier* with a name that identifies the mount point.
2. In the map file, fill in the mount point, options, and location fields as described in [The autofs configuration files](#) section.
3. Register the map file in the master map file, as described in [The autofs configuration files](#) section.
4. Allow the service to re-read the configuration, so it can manage the newly configured **autofs** mount:

```
# systemctl reload autofs.service
```

-
- 5. Try accessing content in the on-demand directory:

```
# ls automounted-directory
```

26.4. AUTOMOUNTING NFS SERVER USER HOME DIRECTORIES WITH AUTOFS SERVICE

This procedure describes how to configure the **autofs** service to mount user home directories automatically.

Prerequisites

- The **autofs** package is installed.
- The **autofs** service is enabled and running.

Procedure

1. Specify the mount point and location of the map file by editing the **/etc/auto.master** file on a server on which you need to mount user home directories. To do so, add the following line into the **/etc/auto.master** file:

```
/home /etc/auto.home
```

2. Create a map file with the name of **/etc/auto.home** on a server on which you need to mount user home directories, and edit the file with the following parameters:

```
* -fstype=nfs,rw,sync host.example.com:/home/&
```

You can skip **fstype** parameter, as it is **nfs** by default. For more information, see **autofs(5)** man page.

3. Reload the **autofs** service:

```
# systemctl reload autofs
```

26.5. OVERRIDING OR AUGMENTING AUTOFS SITE CONFIGURATION FILES

It is sometimes useful to override site defaults for a specific mount point on a client system.

Example 26.3. Initial conditions

For example, consider the following conditions:

- Automounter maps are stored in NIS and the **/etc/nsswitch.conf** file has the following directive:

```
automount: files nis
```

- The **auto.master** file contains:

```
+auto.master
```

- The NIS **auto.master** map file contains:

```
/home auto.home
```

- The NIS **auto.home** map contains:

```
beth  fileserver.example.com:/export/home/beth
joe   fileserver.example.com:/export/home/joe
*     fileserver.example.com:/export/home/&
```

- The **autofs** configuration option **BROWSE_MODE** is set to **yes**:

```
BROWSE_MODE="yes"
```

- The file map **/etc/auto.home** does not exist.

Procedure

This section describes the examples of mounting home directories from a different server and augmenting **auto.home** with only selected entries.

Example 26.4. Mounting home directories from a different server

Given the preceding conditions, let's assume that the client system needs to override the NIS map **auto.home** and mount home directories from a different server.

- In this case, the client needs to use the following **/etc/auto.master** map:

```
/home /etc/auto.home
+auto.master
```

- The **/etc/auto.home** map contains the entry:

```
* host.example.com:/export/home/&
```

Because the automounter only processes the first occurrence of a mount point, the **/home** directory contains the content of **/etc/auto.home** instead of the NIS **auto.home** map.

Example 26.5. Augmenting auto.home with only selected entries

Alternatively, to augment the site-wide **auto.home** map with just a few entries:

1. Create an **/etc/auto.home** file map, and in it put the new entries. At the end, include the NIS **auto.home** map. Then the **/etc/auto.home** file map looks similar to:

```
mydir someserver:/export/mydir
+auto.home
```

2. With these NIS **auto.home** map conditions, listing the content of the **/home** directory outputs:

```
$ ls /home

beth joe mydir
```

This last example works as expected because **autofs** does not include the contents of a file map of the same name as the one it is reading. As such, **autofs** moves on to the next map source in the **nsswitch** configuration.

26.6. USING LDAP TO STORE AUTOMOUNTER MAPS

This procedure configures **autofs** to store automounter maps in LDAP configuration rather than in **autofs** map files.

Prerequisites

- LDAP client libraries must be installed on all systems configured to retrieve automounter maps from LDAP. On Red Hat Enterprise Linux, the **openldap** package should be installed automatically as a dependency of the **autofs** package.

Procedure

1. To configure LDAP access, modify the **/etc/openldap/ldap.conf** file. Ensure that the **BASE**, **URI**, and **schema** options are set appropriately for your site.
2. The most recently established schema for storing automount maps in LDAP is described by the **rfc2307bis** draft. To use this schema, set it in the **/etc/autofs.conf** configuration file by removing the comment characters from the schema definition. For example:

Example 26.6. Setting autofs configuration

```
DEFAULT_MAP_OBJECT_CLASS="automountMap"
DEFAULT_ENTRY_OBJECT_CLASS="automount"
DEFAULT_MAP_ATTRIBUTE="automountMapName"
DEFAULT_ENTRY_ATTRIBUTE="automountKey"
DEFAULT_VALUE_ATTRIBUTE="automountInformation"
```

3. Ensure that all other schema entries are commented in the configuration. The **automountKey** attribute of the **rfc2307bis** schema replaces the **cn** attribute of the **rfc2307** schema. Following is an example of an LDAP Data Interchange Format (LDIF) configuration:

Example 26.7. LDIF Configuration

```
# auto.master, example.com
dn: automountMapName=auto.master,dc=example,dc=com
objectClass: top
objectClass: automountMap
automountMapName: auto.master

# /home, auto.master, example.com
dn: automountMapName=auto.master,dc=example,dc=com
```

```

objectClass: automount
automountKey: /home
automountInformation: auto.home

# auto.home, example.com
dn: automountMapName=auto.home,dc=example,dc=com
objectClass: automountMap
automountMapName: auto.home

# foo, auto.home, example.com
dn: automountKey=foo,automountMapName=auto.home,dc=example,dc=com
objectClass: automount
automountKey: foo
automountInformation: filer.example.com:/export/foo

# /, auto.home, example.com
dn: automountKey=/,automountMapName=auto.home,dc=example,dc=com
objectClass: automount
automountKey: /
automountInformation: filer.example.com:/export/&

```

Additional resources

- [The rfc2307bis draft](#)

26.7. USING SYSTEMD.AUTOMOUNT TO MOUNT A FILE SYSTEM ON DEMAND WITH /ETC/FSTAB

This procedure shows how to mount a file system on demand using the automount systemd units when mount point is defined in **/etc/fstab**. You have to add an automount unit for each mount and enable it.

Procedure

1. Add desired fstab entry as documented in [Persistently mounting file systems](#). For example:

```
/dev/disk/by-id/da875760-edb9-4b82-99dc-5f4b1ff2e5f4 /mount/point xfs defaults 0 0
```

2. Add **x-systemd.automount** to the options field of entry created in the previous step.
3. Load newly created units so that your system registers the new configuration:

```
# systemctl daemon-reload
```

4. Start the automount unit:

```
# systemctl start mount-point.automount
```

Verification

1. Check that **mount-point.automount** is running:

```
# systemctl status mount-point.automount
```

-
- 2. Check that automounted directory has desired content:

```
# ls /mount/point
```

Additional resources

- **systemd.automount(5)** man page.
- **systemd.mount(5)** man page.
- [Introduction to systemd](#).

26.8. USING SYSTEMD.AUTOMOUNT TO MOUNT A FILE SYSTEM ON DEMAND WITH A MOUNT UNIT

This procedure shows how to mount a file system on demand using the automount systemd units when mount point is defined by a mount unit. You have to add an automount unit for each mount and enable it.

Procedure

1. Create a mount unit. For example:

```
mount-point.mount
[Mount]
What=/dev/disk/by-uuid/f5755511-a714-44c1-a123-cfde0e4ac688
Where=/mount/point
Type=xfs
```

2. Create a unit file with the same name as the mount unit, but with extension **.automount**.
3. Open the file and create an **[Automount]** section. Set the **Where=** option to the mount path:

```
[Automount]
Where=/mount/point
[Install]
WantedBy=multi-user.target
```

4. Load newly created units so that your system registers the new configuration:

```
# systemctl daemon-reload
```

5. Enable and start the automount unit instead:

```
# systemctl enable --now mount-point.automount
```

Verification

1. Check that **mount-point.automount** is running:

```
# systemctl status mount-point.automount
```

2. Check that automounted directory has desired content:

```
# ls /mount/point
```

Additional resources

- **systemd.automount(5)** man page.
- **systemd.mount(5)** man page.
- [Introduction to systemd](#).

CHAPTER 27. USING SSSD COMPONENT FROM IDM TO CACHE THE AUTOFS MAPS

The System Security Services Daemon (SSSD) is a system service to access remote service directories and authentication mechanisms. The data caching is useful in case of the slow network connection. To configure the SSSD service to cache the autofs map, follow the procedures below in this section.

27.1. CONFIGURING AUTOFS MANUALLY TO USE IDM SERVER AS AN LDAP SERVER

This procedure shows how to configure **autofs** to use IdM server as an LDAP server.

Procedure

1. Edit the **/etc/autofs.conf** file to specify the schema attributes that **autofs** searches for:

```
#
# Other common LDAP naming
#
map_object_class = "automountMap"
entry_object_class = "automount"
map_attribute = "automountMapName"
entry_attribute = "automountKey"
value_attribute = "automountInformation"
```



NOTE

User can write the attributes in both lower and upper cases in the **/etc/autofs.conf** file.

2. Optionally, specify the LDAP configuration. There are two ways to do this. The simplest is to let the automount service discover the LDAP server and locations on its own:

```
ldap_uri = "ldap:///dc=example,dc=com"
```

This option requires DNS to contain SRV records for the discoverable servers.

Alternatively, explicitly set which LDAP server to use and the base DN for LDAP searches:

```
ldap_uri = "ldap://ipa.example.com"
search_base = "cn=location,cn=automount,dc=example,dc=com"
```

3. Edit the **/etc/autofs_ldap_auth.conf** file so that autofs allows client authentication with the IdM LDAP server.

- Change **authrequired** to yes.
- Set the principal to the Kerberos host principal for the IdM LDAP server, *host/fqdn@REALM*. The principal name is used to connect to the IdM directory as part of GSS client authentication.

```
<autofs_ldap_sasl_conf
```



```

usetls="no"
tlsrequired="no"
authrequired="yes"
authtype="GSSAPI"
clientprinc="host/server.example.com@EXAMPLE.COM"
/>

```

For more information about host principal, see [Using canonicalized DNS host names in IdM](#) .

If necessary, run **klist -k** to get the exact host principal information.

27.2. CONFIGURING SSSD TO CACHE AUTOFS MAPS

The SSSD service can be used to cache **autofs** maps stored on an IdM server without having to configure **autofs** to use the IdM server at all.

Prerequisites

- The **sssd** package is installed.

Procedure

1. Open the SSSD configuration file:

```
# vim /etc/sssd/sssd.conf
```

2. Add the **autofs** service to the list of services handled by SSSD.

```

[sssd]
domains = ldap
services = nss,pam,autofs

```

3. Create a new **[autofs]** section. You can leave this blank, because the default settings for an **autofs** service work with most infrastructures.

```

[nss]

[pam]

[sudo]

[autofs]

[ssh]

[pac]

```

For more information, see the **sssd.conf** man page.

4. Optionally, set a search base for the **autofs** entries. By default, this is the LDAP search base, but a subtree can be specified in the **ldap_autofs_search_base** parameter.

```
[domain/EXAMPLE]
```

```
ldap_search_base = "dc=example,dc=com"
ldap_autofs_search_base = "ou=automount,dc=example,dc=com"
```

5. Restart SSSD service:

```
# systemctl restart sssd.service
```

6. Check the **/etc/nsswitch.conf** file, so that SSSD is listed as a source for automount configuration:

```
automount: sss files
```

7. Restart **autofs** service:

```
# systemctl restart autofs.service
```

8. Test the configuration by listing a user's **/home** directory, assuming there is a master map entry for **/home**:

```
# ls /home/userName
```

If this does not mount the remote file system, check the **/var/log/messages** file for errors. If necessary, increase the debug level in the **/etc/sysconfig/autofs** file by setting the **logging** parameter to **debug**.

CHAPTER 28. SETTING READ-ONLY PERMISSIONS FOR THE ROOT FILE SYSTEM

Sometimes, you need to mount the root file system (/) with read-only permissions. Example use cases include enhancing security or ensuring data integrity after an unexpected system power-off.

28.1. FILES AND DIRECTORIES THAT ALWAYS RETAIN WRITE PERMISSIONS

For the system to function properly, some files and directories need to retain write permissions. When the root file system is mounted in read-only mode, these files are mounted in RAM using the **tmpfs** temporary file system.

The default set of such files and directories is read from the **/etc/rwtab** file. Note that the **readonly-root** package is required to have this file present in your system.

```
dirs /var/cache/man
dirs /var/gdm
<content truncated>

empty /tmp
empty /var/cache/foomatic
<content truncated>

files /etc/adjtime
files /etc/ntp.conf
<content truncated>
```

Entries in the **/etc/rwtab** file follow this format:

```
copy-method path
```

In this syntax:

- Replace *copy-method* with one of the keywords specifying how the file or directory is copied to tmpfs.
- Replace *path* with the path to the file or directory.

The **/etc/rwtab** file recognizes the following ways in which a file or directory can be copied to **tmpfs**:

empty

An empty path is copied to **tmpfs**. For example:

```
empty /tmp
```

dirs

A directory tree is copied to **tmpfs**, empty. For example:

```
dirs /var/run
```

files

A file or a directory tree is copied to **tmpfs** intact. For example:

```
files /etc/resolv.conf
```

The same format applies when adding custom paths to **/etc/rwtab.d/**.

28.2. CONFIGURING THE ROOT FILE SYSTEM TO MOUNT WITH READ-ONLY PERMISSIONS ON BOOT

With this procedure, the root file system is mounted read-only on all following boots.

Procedure

1. In the **/etc/sysconfig/readonly-root** file, set the **READONLY** option to **yes** to mount the file systems as read-only:

```
READONLY=yes
```

2. Add the **ro** option in the root entry (**/**) in the **/etc/fstab** file:

```
/dev/mapper/luks-c376919e... / xfs x-systemd.device-timeout=0,ro 1 1
```

3. Enable the **ro** kernel option:

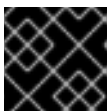
```
# grubby --update-kernel=ALL --args="ro"
```

4. Ensure that the **rw** kernel option is disabled:

```
# grubby --update-kernel=ALL --remove-args="rw"
```

5. If you need to add files and directories to be mounted with write permissions in the **tmpfs** file system, create a text file in the **/etc/rwtab.d/** directory and put the configuration there. For example, to mount the **/etc/example/file** file with write permissions, add this line to the **/etc/rwtab.d/example** file:

```
files /etc/example/file
```



IMPORTANT

Changes made to files and directories in **tmpfs** do not persist across boots.

6. Reboot the system to apply the changes.

Troubleshooting

- If you mount the root file system with read-only permissions by mistake, you can remount it with read-and-write permissions again using the following command:

```
# mount -o remount,rw /
```

CHAPTER 29. LIMITING STORAGE SPACE USAGE ON XFS WITH QUOTAS

You can restrict the amount of disk space available to users or groups by implementing disk quotas. You can also define a warning level at which system administrators are informed before a user consumes too much disk space or a partition becomes full.

The XFS quota subsystem manages limits on disk space (blocks) and file (inode) usage. XFS quotas control or report on usage of these items on a user, group, or directory or project level. Group and project quotas are only mutually exclusive on older non-default XFS disk formats.

When managing on a per-directory or per-project basis, XFS manages the disk usage of directory hierarchies associated with a specific project.

29.1. DISK QUOTAS

In most computing environments, disk space is not infinite. The quota subsystem provides a mechanism to control usage of disk space.

You can configure disk quotas for individual users as well as user groups on the local file systems. This makes it possible to manage the space allocated for user-specific files (such as email) separately from the space allocated to the projects that a user works on. The quota subsystem warns users when they exceed their allotted limit, but allows some extra space for current work (hard limit/soft limit).

If quotas are implemented, you need to check if the quotas are exceeded and make sure the quotas are accurate. If users repeatedly exceed their quotas or consistently reach their soft limits, a system administrator can either help the user determine how to use less disk space or increase the user's disk quota.

You can set quotas to control:

- The number of consumed disk blocks.
- The number of inodes, which are data structures that contain information about files in UNIX file systems. Because inodes store file-related information, this allows control over the number of files that can be created.

29.2. THE `xfs_quota` TOOL

You can use the **`xfs_quota`** tool to manage quotas on XFS file systems. In addition, you can use XFS file systems with limit enforcement turned off as an effective disk usage accounting system.

The XFS quota system differs from other file systems in a number of ways. Most importantly, XFS considers quota information as file system metadata and uses journaling to provide a higher level guarantee of consistency.

Additional resources

- **`xfs_quota(8)`** man page.

29.3. FILE SYSTEM QUOTA MANAGEMENT IN XFS

The XFS quota subsystem manages limits on disk space (blocks) and file (inode) usage. XFS quotas control or report on usage of these items on a user, group, or directory or project level. Group and project quotas are only mutually exclusive on older non-default XFS disk formats.

When managing on a per-directory or per-project basis, XFS manages the disk usage of directory hierarchies associated with a specific project.

29.4. ENABLING DISK QUOTAS FOR XFS

This procedure enables disk quotas for users, groups, and projects on an XFS file system. Once quotas are enabled, the **xfs_quota** tool can be used to set limits and report on disk usage.

Procedure

1. Enable quotas for users:

```
# mount -o uquota /dev/xvdb1 /xfs
```

Replace **uquota** with **uqnoenforce** to allow usage reporting without enforcing any limits.

2. Enable quotas for groups:

```
# mount -o gquota /dev/xvdb1 /xfs
```

Replace **gquota** with **gqnoenforce** to allow usage reporting without enforcing any limits.

3. Enable quotas for projects:

```
# mount -o pquota /dev/xvdb1 /xfs
```

Replace **pquota** with **pqnoenforce** to allow usage reporting without enforcing any limits.

4. Alternatively, include the quota mount options in the **/etc/fstab** file. The following example shows entries in the **/etc/fstab** file to enable quotas for users, groups, and projects, respectively, on an XFS file system. These examples also mount the file system with read/write permissions:

```
# vim /etc/fstab
/dev/xvdb1 /xfs xfs rw,quota 0 0
/dev/xvdb1 /xfs xfs rw,gquota 0 0
/dev/xvdb1 /xfs xfs rw,prjquota 0 0
```

Additional resources

- **mount(8)** man page.
- **xfs_quota(8)** man page.

29.5. REPORTING XFS USAGE

You can use the **xfs_quota** tool to set limits and report on disk usage. By default, **xfs_quota** is run interactively, and in basic mode. Basic mode subcommands simply report usage, and are available to all users.

Prerequisites

- Quotas have been enabled for the XFS file system. See [Enabling disk quotas for XFS](#).

Procedure

1. Start the **xfs_quota** shell:

```
# xfs_quota
```

2. Show usage and limits for the given user:

```
# xfs_quota> quota username
```

3. Show free and used counts for blocks and inodes:

```
# xfs_quota> df
```

4. Run the help command to display the basic commands available with **xfs_quota**.

```
# xfs_quota> help
```

5. Specify **q** to exit **xfs_quota**.

```
# xfs_quota> q
```

Additional resources

- **xfs_quota(8)** man page.

29.6. MODIFYING XFS QUOTA LIMITS

Start the **xfs_quota** tool with the **-x** option to enable expert mode and run the administrator commands, which allow modifications to the quota system. The subcommands of this mode allow actual configuration of limits, and are available only to users with elevated privileges.

Prerequisites

- Quotas have been enabled for the XFS file system. See [Enabling disk quotas for XFS](#).

Procedure

1. Start the **xfs_quota** shell with the **-x** option to enable expert mode:

```
# xfs_quota -x
```

2. Report quota information for a specific file system:

```
# xfs_quota> report /path
```

For example, to display a sample quota report for **/home** (on **/dev/blockdevice**), use the command **report -h /home**. This displays output similar to the following:

```
User quota on /home (/dev/blockdevice)
Blocks
User ID    Used  Soft  Hard Warn/Grace
-----
root       0    0    0 00 [-----]
testuser  103.4G    0    0 00 [-----]
```

3. Modify quota limits:

```
# xfs_quota> limit isoft=500m ihard=700m user /path
```

For example, to set a soft and hard inode count limit of 500 and 700 respectively for user **john**, whose home directory is **/home/john**, use the following command:

```
# xfs_quota -x -c 'limit isoft=500 ihard=700 john' /home/
```

In this case, pass **mount_point** which is the mounted xfs file system.

4. Run the help command to display the expert commands available with **xfs_quota -x**:

```
# xfs_quota> help
```

Additional resources

- **xfs_quota(8)** man page.

29.7. SETTING PROJECT LIMITS FOR XFS

This procedure configures limits for project-controlled directories.

Procedure

1. Add the project-controlled directories to **/etc/projects**. For example, the following adds the **/var/log** path with a unique ID of 11 to **/etc/projects**. Your project ID can be any numerical value mapped to your project.

```
# echo 11:/var/log >> /etc/projects
```

2. Add project names to **/etc/projid** to map project IDs to project names. For example, the following associates a project called **logfiles** with the project ID of 11 as defined in the previous step.

```
# echo logfiles:11 >> /etc/projid
```

3. Initialize the project directory. For example, the following initializes the project directory **/var**:

```
# xfs_quota -x -c 'project -s logfiles' /var
```

4. Configure quotas for projects with initialized directories:

```
# xfs_quota -x -c 'limit -p bhard=1g logfiles' /var
```


Additional resources

- **xfstool(8)** man page.
- **projid(5)** man page.
- **projects(5)** man page.

CHAPTER 30. LIMITING STORAGE SPACE USAGE ON EXT4 WITH QUOTAS

You have to enable disk quotas on your system before you can assign them. You can assign disk quotas per user, per group or per project. However, if there is a soft limit set, you can exceed these quotas for a configurable period of time, known as the grace period.

30.1. INSTALLING THE QUOTA TOOL

You must install the **quota** RPM package to implement disk quotas.

Procedure

- Install the **quota** package:

```
# dnf install quota
```

30.2. ENABLING QUOTA FEATURE ON FILE SYSTEM CREATION

This procedure describes how to enable quotas on file system creation.

Procedure

1. Enable quotas on file system creation:

```
# mkfs.ext4 -O quota /dev/sda
```



NOTE

Only user and group quotas are enabled and initialized by default.

2. Change the defaults on file system creation:

```
# mkfs.ext4 -O quota -E quotatype=usrquota:grpquota:prjquota /dev/sda
```

3. Mount the file system:

```
# mount /dev/sda
```

Additional resources

- **ext4(5)** man page.

30.3. ENABLING QUOTA FEATURE ON EXISTING FILE SYSTEMS

This procedure describes how to enable the quota feature on existing file system using the **tune2fs** command.

Procedure

1. Unmount the file system:

```
# umount /dev/sda
```

2. Enable quotas on existing file system:

```
# tune2fs -O quota /dev/sda
```



NOTE

Only user and group quotas are initialized by default.

3. Change the defaults:

```
# tune2fs -Q usrquota,grpquota,prjquota /dev/sda
```

4. Mount the file system:

```
# mount /dev/sda
```

Additional resources

- **ext4(5)** man page.

30.4. ENABLING QUOTA ENFORCEMENT

The quota accounting is enabled by default after mounting the file system without any additional options, but quota enforcement is not.

Prerequisites

- Quota feature is enabled and the default quotas are initialized.

Procedure

- Enable quota enforcement by **quotaon** for the user quota:

```
# mount /dev/sda /mnt
```

```
# quotaon /mnt
```



NOTE

The quota enforcement can be enabled at mount time using **usrquota**, **grpquota**, or **prjquota** mount options.

```
# mount -o usrquota,grpquota,prjquota /dev/sda /mnt
```

- Enable user, group, and project quotas for all file systems:

```
# quotaon -vaugP
```

- If neither of the **-u**, **-g**, or **-P** options are specified, only the user quotas are enabled.
- If only **-g** option is specified, only group quotas are enabled.
- If only **-P** option is specified, only project quotas are enabled.
- Enable quotas for a specific file system, such as **/home**:

```
# quotaon -vugP /home
```

Additional resources

- **quotaon(8)** man page.

30.5. ASSIGNING QUOTAS PER USER

The disk quotas are assigned to users with the **edquota** command.



NOTE

The text editor defined by the **EDITOR** environment variable is used by **edquota**. To change the editor, set the **EDITOR** environment variable in your **~/.bash_profile** file to the full path of the editor of your choice.

Prerequisites

- User must exist prior to setting the user quota.

Procedure

1. Assign the quota for a user:

```
# edquota username
```

Replace *username* with the user to which you want to assign the quotas.

For example, if you enable a quota for the **/dev/sda** partition and execute the command **edquota testuser**, the following is displayed in the default editor configured on the system:

```
Disk quotas for user testuser (uid 501):
Filesystem blocks soft hard inodes soft hard
/dev/sda 44043 0 0 37418 0 0
```

2. Change the desired limits.
If any of the values are set to 0, limit is not set. Change them in the text editor.

For example, the following shows the soft and hard block limits for the testuser have been set to 50000 and 55000 respectively.

Disk quotas for user `testuser` (uid 501):

```
Filesystem blocks soft hard inodes soft hard
/dev/sda 44043 50000 55000 37418 0 0
```

- The first column is the name of the file system that has a quota enabled for it.
- The second column shows how many blocks the user is currently using.
- The next two columns are used to set soft and hard block limits for the user on the file system.
- The **inodes** column shows how many inodes the user is currently using.
- The last two columns are used to set the soft and hard inode limits for the user on the file system.
 - The hard block limit is the absolute maximum amount of disk space that a user or group can use. Once this limit is reached, no further disk space can be used.
 - The soft block limit defines the maximum amount of disk space that can be used. However, unlike the hard limit, the soft limit can be exceeded for a certain amount of time. That time is known as the *grace period*. The grace period can be expressed in seconds, minutes, hours, days, weeks, or months.

Verification steps

- Verify that the quota for the user has been set:

```
# quota -v testuser
```

Disk quotas for user `testuser`:

```
Filesystem blocks quota limit grace files quota limit grace
/dev/sda 1000* 1000 1000 0 0 0
```

30.6. ASSIGNING QUOTAS PER GROUP

You can assign quotas on a per-group basis.

Prerequisites

- Group must exist prior to setting the group quota.

Procedure

1. Set a group quota:

```
# edquota -g groupname
```

For example, to set a group quota for the **devel** group:

```
# edquota -g devel
```

This command displays the existing quota for the group in the text editor:

```
Disk quotas for group devel (gid 505):
Filesystem blocks soft hard inodes soft hard
/dev/sda 440400 0 0 37418 0 0
```

2. Modify the limits and save the file.

Verification steps

- Verify that the group quota is set:

```
# quota -vg groupname
```

30.7. ASSIGNING QUOTAS PER PROJECT

This procedure assigns quotas per project.

Prerequisites

- Project quota is enabled on your file system.

Procedure

1. Add the project-controlled directories to **/etc/projects**. For example, the following adds the **/var/log** path with a unique ID of 11 to **/etc/projects**. Your project ID can be any numerical value mapped to your project.

```
# echo 11:/var/log >> /etc/projects
```

2. Add project names to **/etc/projid** to map project IDs to project names. For example, the following associates a project called **Logs** with the project ID of 11 as defined in the previous step.

```
# echo Logs:11 >> /etc/projid
```

3. Set the desired limits:

```
# edquota -P 11
```



NOTE

You can choose the project either by its project ID (**11** in this case), or by its name (**Logs** in this case).

4. Using **quotaon**, enable quota enforcement:
See [Enabling quota enforcement](#).

Verification steps

- Verify that the project quota is set:

```
# quota -vP 11
```

**NOTE**

You can verify either by the project ID, or by the project name.

Additional resources

- **edquota(8)** man page.
- **projid(5)** man page.
- **projects(5)** man page.

30.8. SETTING THE GRACE PERIOD FOR SOFT LIMITS

If a given quota has soft limits, you can edit the grace period, which is the amount of time for which a soft limit can be exceeded. You can set the grace period for users, groups, or projects.

Procedure

- Edit the grace period:

```
# edquota -t
```

**IMPORTANT**

While other **edquota** commands operate on quotas for a particular user, group, or project, the **-t** option operates on every file system with quotas enabled.

Additional resources

- **edquota(8)** man page.

30.9. TURNING FILE SYSTEM QUOTAS OFF

Use **quotaoff** to turn disk quota enforcement off on the specified file systems. Quota accounting stays enabled after executing this command.

Procedure

- To turn all user and group quotas off:

```
# quotaoff -vaugP
```

- If neither of the **-u**, **-g**, or **-P** options are specified, only the user quotas are disabled.
- If only **-g** option is specified, only group quotas are disabled.
- If only **-P** option is specified, only project quotas are disabled.
- The **-v** switch causes verbose status information to display as the command executes.

Additional resources

- **quotaoff(8)** man page.

30.10. REPORTING ON DISK QUOTAS

You can create a disk quota report using the **repquota** utility.

Procedure

1. Run the **repquota** command:

```
# repquota
```

For example, the command **repquota /dev/sda** produces this output:

```
*** Report for user quotas on device /dev/sda
Block grace time: 7days; Inode grace time: 7days
Block limits  File limits
User  used soft hard grace used soft hard grace
-----
root  --   36   0   0         4   0   0
kristin --  540   0   0        125   0   0
testuser -- 440400 500000 550000      37418   0   0
```

2. View the disk usage report for all quota-enabled file systems:

```
# repquota -augP
```

The **--** symbol displayed after each user determines whether the block or inode limits have been exceeded. If either soft limit is exceeded, a **+** character appears in place of the corresponding **-** character. The first **-** character represents the block limit, and the second represents the inode limit.

The **grace** columns are normally blank. If a soft limit has been exceeded, the column contains a time specification equal to the amount of time remaining on the grace period. If the grace period has expired, **none** appears in its place.

Additional resources

The **repquota(8)** man page for more information.

CHAPTER 31. DISCARDING UNUSED BLOCKS

You can perform or schedule discard operations on block devices that support them. The block discard operation communicates to the underlying storage which filesystem blocks are no longer in use by the mounted filesystem. Block discard operations allow SSDs to optimize garbage collection routines, and they can inform thinly-provisioned storage to repurpose unused physical blocks.

Requirements

- The block device underlying the file system must support physical discard operations. Physical discard operations are supported if the value in the `/sys/block/<device>/queue/discard_max_bytes` file is not zero.

31.1. TYPES OF BLOCK DISCARD OPERATIONS

You can run discard operations using different methods:

Batch discard

Is triggered explicitly by the user and discards all unused blocks in the selected file systems.

Online discard

Is specified at mount time and triggers in real time without user intervention. Online discard operations discard only blocks that are transitioning from the **used** to the **free** state.

Periodic discard

Are batch operations that are run regularly by a **systemd** service.

All types are supported by the XFS and ext4 file systems.

Recommendations

Red Hat recommends that you use batch or periodic discard.

Use online discard only if:

- the system's workload is such that batch discard is not feasible, or
- online discard operations are necessary to maintain performance.

31.2. PERFORMING BATCH BLOCK DISCARD

You can perform a batch block discard operation to discard unused blocks on a mounted file system.

Prerequisites

- The file system is mounted.
- The block device underlying the file system supports physical discard operations.

Procedure

- Use the **fstrim** utility:
 - To perform discard only on a selected file system, use:

```
# fstrim mount-point
```

- To perform discard on all mounted file systems, use:

```
# fstrim --all
```

If you execute the **fstrim** command on:

- a device that does not support discard operations, or
- a logical device (LVM or MD) composed of multiple devices, where any one of the device does not support discard operations,

the following message displays:

```
# fstrim /mnt/non_discard
```

```
fstrim: /mnt/non_discard: the discard operation is not supported
```

Additional resources

- **fstrim(8)** man page.

31.3. ENABLING ONLINE BLOCK DISCARD

You can perform online block discard operations to automatically discard unused blocks on all supported file systems.

Procedure

- Enable online discard at mount time:
 - When mounting a file system manually, add the **-o discard** mount option:

```
# mount -o discard device mount-point
```
 - When mounting a file system persistently, add the **discard** option to the mount entry in the **/etc/fstab** file.

Additional resources

- **mount(8)** man page.
- **fstab(5)** man page.

31.4. ENABLING PERIODIC BLOCK DISCARD

You can enable a **systemd** timer to regularly discard unused blocks on all supported file systems.

Procedure

- Enable and start the **systemd** timer:

```
# systemctl enable --now fstrim.timer
Created symlink /etc/systemd/system/timers.target.wants/fstrim.timer →
/usr/lib/systemd/system/fstrim.timer.
```

Verification

- Verify the status of the timer:

```
# systemctl status fstrim.timer
fstrim.timer - Discard unused blocks once a week
  Loaded: loaded (/usr/lib/systemd/system/fstrim.timer; enabled; vendor preset: disabled)
  Active: active (waiting) since Wed 2023-05-17 13:24:41 CEST; 3min 15s ago
  Trigger: Mon 2023-05-22 01:20:46 CEST; 4 days left
  Docs: man:fstrim

May 17 13:24:41 localhost.localdomain systemd[1]: Started Discard unused blocks once a
week.
```

CHAPTER 32. SETTING UP STRATIS FILE SYSTEMS

Stratis runs as a service to manage pools of physical storage devices, simplifying local storage management with ease of use while helping you set up and manage complex storage configurations.



IMPORTANT

Stratis is a Technology Preview feature only. Technology Preview features are not supported with Red Hat production service level agreements (SLAs) and might not be functionally complete. Red Hat does not recommend using them in production. These features provide early access to upcoming product features, enabling customers to test functionality and provide feedback during the development process. For more information about the support scope of Red Hat Technology Preview features, see <https://access.redhat.com/support/offerings/techpreview>.

32.1. WHAT IS STRATIS

Stratis is a local storage-management solution for Linux. It is focused on simplicity and ease of use, and gives you access to advanced storage features.

Stratis makes the following activities easier:

- Initial configuration of storage
- Making changes later
- Using advanced storage features

Stratis is a local storage management system that supports advanced storage features. The central concept of Stratis is a storage *pool*. This pool is created from one or more local disks or partitions, and file systems are created from the pool.

The pool enables many useful features, such as:

- File system snapshots
- Thin provisioning
- Tiering
- Encryption

Additional resources

- [Stratis website](#)

32.2. COMPONENTS OF A STRATIS VOLUME

Learn about the components that comprise a Stratis volume.

Externally, Stratis presents the following volume components in the command-line interface and the API:

blockdev

Block devices, such as a disk or a disk partition.

pool

Composed of one or more block devices.

A pool has a fixed total size, equal to the size of the block devices.

The pool contains most Stratis layers, such as the non-volatile data cache using the **dm-cache** target.

Stratis creates a **/dev/stratis/my-pool/** directory for each pool. This directory contains links to devices that represent Stratis file systems in the pool.

filesystem

Each pool can contain one or more file systems, which store files.

File systems are thinly provisioned and do not have a fixed total size. The actual size of a file system grows with the data stored on it. If the size of the data approaches the virtual size of the file system, Stratis grows the thin volume and the file system automatically.

The file systems are formatted with XFS.



IMPORTANT

Stratis tracks information about file systems created using Stratis that XFS is not aware of, and changes made using XFS do not automatically create updates in Stratis. Users must not reformat or reconfigure XFS file systems that are managed by Stratis.

Stratis creates links to file systems at the **/dev/stratis/my-pool/my-fs** path.



NOTE

Stratis uses many Device Mapper devices, which show up in **dmsetup** listings and the **/proc/partitions** file. Similarly, the **lsblk** command output reflects the internal workings and layers of Stratis.

32.3. BLOCK DEVICES USABLE WITH STRATIS

Storage devices that can be used with Stratis.

Supported devices

Stratis pools have been tested to work on these types of block devices:

- LUKS
- LVM logical volumes
- MD RAID
- DM Multipath
- iSCSI
- HDDs and SSDs
- NVMe devices

Unsupported devices

Because Stratis contains a thin-provisioning layer, Red Hat does not recommend placing a Stratis pool on block devices that are already thinly-provisioned.

32.4. INSTALLING STRATIS

Install the required packages for Stratis.

Procedure

1. Install packages that provide the Stratis service and command-line utilities:

```
# dnf install stratisd stratis-cli
```

2. Verify that the **stratisd** service is enabled:

```
# systemctl enable --now stratisd
```

32.5. CREATING AN UNENCRYPTED STRATIS POOL

You can create an unencrypted Stratis pool from one or more block devices.

Prerequisites

- Stratis is installed. For more information, see [Installing Stratis](#).
- The **stratisd** service is running.
- The block devices on which you are creating a Stratis pool are not in use and are not mounted.
- Each block device on which you are creating a Stratis pool is at least 1 GB.
- On the IBM Z architecture, the **/dev/dasd*** block devices must be partitioned. Use the partition device for creating the Stratis pool.

For information about partitioning DASD devices, see [Configuring a Linux instance on IBM Z](#).



NOTE

You cannot encrypt an unencrypted Stratis pool.

Procedure

1. Erase any file system, partition table, or RAID signatures that exist on each block device that you want to use in the Stratis pool:

```
# wipefs --all block-device
```

where ***block-device*** is the path to the block device; for example, **/dev/sdb**.

2. Create the new unencrypted Stratis pool on the selected block device:

```
# stratis pool create my-pool block-device
```

where **block-device** is the path to an empty or wiped block device.



NOTE

Specify multiple block devices on a single line:

```
# stratis pool create my-pool block-device-1 block-device-2
```

3. Verify that the new Stratis pool was created:

```
# stratis pool list
```

32.6. CREATING AN ENCRYPTED STRATIS POOL

To secure your data, you can create an encrypted Stratis pool from one or more block devices.

When you create an encrypted Stratis pool, the kernel keyring is used as the primary encryption mechanism. After subsequent system reboots this kernel keyring is used to unlock the encrypted Stratis pool.

When creating an encrypted Stratis pool from one or more block devices, note the following:

- Each block device is encrypted using the **cryptsetup** library and implements the **LUKS2** format.
- Each Stratis pool can either have a unique key or share the same key with other pools. These keys are stored in the kernel keyring.
- The block devices that comprise a Stratis pool must be either all encrypted or all unencrypted. It is not possible to have both encrypted and unencrypted block devices in the same Stratis pool.
- Block devices added to the data tier of an encrypted Stratis pool are automatically encrypted.

Prerequisites

- Stratis v2.1.0 or later is installed. For more information, see [Installing Stratis](#).
- The **stratisd** service is running.
- The block devices on which you are creating a Stratis pool are not in use and are not mounted.
- The block devices on which you are creating a Stratis pool are at least 1GB in size each.
- On the IBM Z architecture, the **/dev/dasd*** block devices must be partitioned. Use the partition in the Stratis pool.

For information about partitioning DASD devices, see [Configuring a Linux instance on IBM Z](#).

Procedure

1. Erase any file system, partition table, or RAID signatures that exist on each block device that you want to use in the Stratis pool:

```
# wipefs --all block-device
```

where **block-device** is the path to the block device; for example, **/dev/sdb**.

- If you have not created a key set already, run the following command and follow the prompts to create a key set to use for the encryption.

```
# stratis key set --capture-key key-description
```

where **key-description** is a reference to the key that gets created in the kernel keyring.

- Create the encrypted Stratis pool and specify the key description to use for the encryption. You can also specify the key path using the **--keyfile-path** option instead of using the **key-description** option.

```
# stratis pool create --key-desc key-description my-pool block-device
```

where

key-description

References the key that exists in the kernel keyring, which you created in the previous step.

my-pool

Specifies the name of the new Stratis pool.

block-device

Specifies the path to an empty or wiped block device.



NOTE

Specify multiple block devices on a single line:

```
# stratis pool create --key-desc key-description my-pool block-device-1  
block-device-2
```

- Verify that the new Stratis pool was created:

```
# stratis pool list
```

32.7. SETTING OVERPROVISIONING MODE IN STRATIS FILESYSTEM

A storage stack can reach a state of overprovision. If the file system size becomes bigger than the pool backing it, the pool becomes full. To prevent this, disable overprovisioning, which ensures that the size of all filesystems on the pool does not exceed the available physical storage provided by the pool. If you use Stratis for critical applications or the root filesystem, this mode prevents certain failure cases.

If you enable overprovisioning, an API signal notifies you when your storage has been fully allocated. The notification serves as a warning to the user to inform them that when all the remaining pool space fills up, Stratis has no space left to extend to.

Prerequisites

- Stratis is installed. For more information, see [Installing Stratis](#).

Procedure

To set up the pool correctly, you have two possibilities:

1. Create a pool from one or more block devices:

```
# stratis pool create --no-overprovision pool-name /dev/sdb
```

- By using the **--no-overprovision** option, the pool cannot allocate more logical space than actual available physical space.

2. Set overprovisioning mode in the existing pool:

```
# stratis pool overprovision pool-name <yes|no>
```

- If set to "yes", you enable overprovisioning to the pool. This means that the sum of the logical sizes of the Stratis filesystems, supported by the pool, can exceed the amount of available data space.

Verification

1. Run the following to view the full list of Stratis pools:

```
# stratis pool list
```

Name	Total Physical	Properties	UUID	Alerts
<i>pool-name</i>	1.42 TiB / 23.96 MiB / 1.42 TiB	~Ca,~Cr,~Op	cb7cb4d8-9322-4ac4-a6fd-eb7ae9e1e540	

2. Check if there is an indication of the pool overprovisioning mode flag in the **stratis pool list** output. The " ~ " is a math symbol for "NOT", so **~Op** means no-overprovisioning.
3. Optional: Run the following to check overprovisioning on a specific pool:

```
# stratis pool overprovision pool-name yes
```

```
# stratis pool list
```

Name	Total Physical	Properties	UUID	Alerts
<i>pool-name</i>	1.42 TiB / 23.96 MiB / 1.42 TiB	~Ca,~Cr,~Op	cb7cb4d8-9322-4ac4-a6fd-eb7ae9e1e540	

Additional resources

- [The Stratis Storage webpage](#).

32.8. BINDING A STRATIS POOL TO NBDE

Binding an encrypted Stratis pool to Network Bound Disk Encryption (NBDE) requires a Tang server. When a system containing the Stratis pool reboots, it connects with the Tang server to automatically unlock the encrypted pool without you having to provide the kernel keyring description.

**NOTE**

Binding a Stratis pool to a supplementary Clevis encryption mechanism does not remove the primary kernel keyring encryption.

Prerequisites

- Stratis v2.3.0 or later is installed. For more information, see [Installing Stratis](#).
- The **stratisd** service is running.
- You have created an encrypted Stratis pool, and you have the key description of the key that was used for the encryption. For more information, see [Creating an encrypted Stratis pool](#).
- You can connect to the Tang server. For more information, see [Deploying a Tang server with SELinux in enforcing mode](#)

Procedure

- Bind an encrypted Stratis pool to NBDE:

```
# stratis pool bind nbde --trust-url my-pool tang-server
```

where

my-pool

Specifies the name of the encrypted Stratis pool.

tang-server

Specifies the IP address or URL of the Tang server.

Additional resources

- [Configuring automated unlocking of encrypted volumes using policy-based decryption](#)

32.9. BINDING A STRATIS POOL TO TPM

When you bind an encrypted Stratis pool to the Trusted Platform Module (TPM) 2.0, the system containing the pool reboots, and the pool is automatically unlocked without you having to provide the kernel keyring description.

Prerequisites

- Stratis v2.3.0 or later is installed. For more information, see [Installing Stratis](#).
- The **stratisd** service is running.
- You have created an encrypted Stratis pool. For more information, see [Creating an encrypted Stratis pool](#).

Procedure

- Bind an encrypted Stratis pool to TPM:

```
# stratis pool bind tpm my-pool key-description
```

■

where

my-pool

Specifies the name of the encrypted Stratis pool.

key-description

References the key that exists in the kernel keyring, which was generated when you created the encrypted Stratis pool.

32.10. UNLOCKING AN ENCRYPTED STRATIS POOL WITH KERNEL KEYRING

After a system reboot, your encrypted Stratis pool or the block devices that comprise it might not be visible. You can unlock the pool using the kernel keyring that was used to encrypt the pool.

Prerequisites

- Stratis v2.1.0 is installed. For more information, see [Installing Stratis](#).
- The **stratisd** service is running.
- You have created an encrypted Stratis pool. For more information, see [Creating an encrypted Stratis pool](#).

Procedure

1. Re-create the key set using the same key description that was used previously:

```
# stratis key set --capture-key key-description
```

where *key-description* references the key that exists in the kernel keyring, which was generated when you created the encrypted Stratis pool.

2. Verify that the Stratis pool is visible:

```
# stratis pool list
```

32.11. UNBINDING A STRATIS POOL FROM SUPPLEMENTARY ENCRYPTION

When you unbind an encrypted Stratis pool from a supported supplementary encryption mechanism, the primary kernel keyring encryption remains in place. This is not true for pools that are created with Clevis encryption from the start.

Prerequisites

- Stratis v2.3.0 or later is installed on your system. For more information, see [Installing Stratis](#).
- You have created an encrypted Stratis pool. For more information, see [Creating an encrypted Stratis pool](#).
- The encrypted Stratis pool is bound to a supported supplementary encryption mechanism.

Procedure

- Unbind an encrypted Stratis pool from a supplementary encryption mechanism:

```
# stratis pool unbind clevis my-pool
```

where

my-pool specifies the name of the Stratis pool you want to unbind.

Additional resources

- [Binding an encrypted Stratis pool to NBDE](#)
- [Binding an encrypted Stratis pool to TPM](#)

32.12. STARTING AND STOPPING STRATIS POOL

You can start and stop Stratis pools. This gives you the option to disassemble or bring down all the objects that were used to construct the pool, such as filesystems, cache devices, thin pool, and encrypted devices. Note that if the pool actively uses any device or filesystem, it might issue a warning and not be able to stop.

The stopped state is recorded in the pool's metadata. These pools do not start on the following boot, until the pool receives a start command.

Prerequisites

- Stratis is installed. For more information, see [Installing Stratis](#).
- The **stratisd** service is running.
- You have created either an unencrypted or an encrypted Stratis pool. See [Creating an unencrypted Stratis pool](#)

or [Creating an encrypted Stratis pool](#).

Procedure

- Use the following command to start the Stratis pool. The **--unlock-method** option specifies the method of unlocking the pool if it is encrypted:

```
# stratis pool start pool-uuid --unlock-method <keyring|clevis>
```

- Alternatively, use the following command to stop the Stratis pool. This tears down the storage stack but leaves all metadata intact:

```
# stratis pool stop pool-name
```

Verification steps

- Use the following command to list all pools on the system:

```
# stratis pool list
```

- Use the following command to list all not previously started pools. If the UUID is specified, the command prints detailed information about the pool corresponding to the UUID:

```
# stratis pool list --stopped --uuid UUID
```

32.13. CREATING A STRATIS FILE SYSTEM

Create a Stratis file system on an existing Stratis pool.

Prerequisites

- Stratis is installed. For more information, see [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis pool. See [Creating an unencrypted Stratis pool](#)

or [Creating an encrypted Stratis pool](#).

Procedure

1. To create a Stratis file system on a pool, use:

```
# stratis filesystem create --size number-and-unit my-pool my-fs
```

where

number-and-unit

Specifies the size of a file system. The specification format must follow the standard size specification format for input, that is B, KiB, MiB, GiB, TiB or PiB.

my-pool

Specifies the name of the Stratis pool.

my-fs

Specifies an arbitrary name for the file system.

For example:

Example 32.1. Creating a Stratis file system

```
# stratis filesystem create --size 10GiB pool1 filesystem1
```

Verification steps

- List file systems within the pool to check if the Stratis filesystem is created:

```
# stratis fs list my-pool
```

Additional resources

- [Mounting a Stratis file system](#).

32.14. MOUNTING A STRATIS FILE SYSTEM

Mount an existing Stratis file system to access the content.

Prerequisites

- Stratis is installed. For more information, see [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis file system. For more information, see [Creating a Stratis filesystem](#).

Procedure

- To mount the file system, use the entries that Stratis maintains in the **/dev/stratis/** directory:

```
# mount /dev/stratis/my-pool/my-fs mount-point
```

The file system is now mounted on the *mount-point* directory and ready to use.

Additional resources

- [Creating a Stratis file system](#).

32.15. PERSISTENTLY MOUNTING A STRATIS FILE SYSTEM

This procedure persistently mounts a Stratis file system so that it is available automatically after booting the system.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis file system. See [Creating a Stratis filesystem](#).

Procedure

1. Determine the UUID attribute of the file system:

```
$ lsblk --output=UUID /dev/stratis/my-pool/my-fs
```

For example:

Example 32.2. Viewing the UUID of Stratis file system

```
$ lsblk --output=UUID /dev/stratis/my-pool/fs1
```

```
UUID
a1f0b64a-4ebb-4d4e-9543-b1d79f600283
```

2. If the mount point directory does not exist, create it:

```
# mkdir --parents mount-point
```

3. As root, edit the **/etc/fstab** file and add a line for the file system, identified by the UUID. Use **xfs** as the file system type and add the **x-systemd.requires=stratisd.service** option. For example:

Example 32.3. The /fs1 mount point in /etc/fstab

```
UUID=a1f0b64a-4ebb-4d4e-9543-b1d79f600283 /fs1 xfs defaults,x-  
systemd.requires=stratisd.service 0 0
```

4. Regenerate mount units so that your system registers the new configuration:

```
# systemctl daemon-reload
```

5. Try mounting the file system to verify that the configuration works:

```
# mount mount-point
```

Additional resources

- [Persistently mounting file systems](#)

32.16. SETTING UP NON-ROOT STRATIS FILESYSTEMS IN /ETC/FSTAB USING A SYSTEMD SERVICE

You can manage setting up non-root filesystems in **/etc/fstab** using a systemd service.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis file system. See [Creating a Stratis filesystem](#).

Procedure

- For all non-root Stratis filesystems, use:

```
# /dev/stratis/[STRATIS_SYMLINK] [MOUNT_POINT] xfs defaults, x-  
systemd.requires=stratis-fstab-setup@[POOL_UUID].service,x-systemd.after=stratis-stab-  
setup@[POOL_UUID].service <dump_value> <fsck_value>
```

Additional resources

- [Persistently mounting file systems](#).

CHAPTER 33. EXTENDING A STRATIS VOLUME WITH ADDITIONAL BLOCK DEVICES

You can attach additional block devices to a Stratis pool to provide more storage capacity for Stratis file systems.



IMPORTANT

Stratis is a Technology Preview feature only. Technology Preview features are not supported with Red Hat production service level agreements (SLAs) and might not be functionally complete. Red Hat does not recommend using them in production. These features provide early access to upcoming product features, enabling customers to test functionality and provide feedback during the development process. For more information about the support scope of Red Hat Technology Preview features, see <https://access.redhat.com/support/offerings/techpreview>.

33.1. COMPONENTS OF A STRATIS VOLUME

Learn about the components that comprise a Stratis volume.

Externally, Stratis presents the following volume components in the command-line interface and the API:

blockdev

Block devices, such as a disk or a disk partition.

pool

Composed of one or more block devices.

A pool has a fixed total size, equal to the size of the block devices.

The pool contains most Stratis layers, such as the non-volatile data cache using the **dm-cache** target.

Stratis creates a **/dev/stratis/my-pool/** directory for each pool. This directory contains links to devices that represent Stratis file systems in the pool.

filesystem

Each pool can contain one or more file systems, which store files.

File systems are thinly provisioned and do not have a fixed total size. The actual size of a file system grows with the data stored on it. If the size of the data approaches the virtual size of the file system, Stratis grows the thin volume and the file system automatically.

The file systems are formatted with XFS.



IMPORTANT

Stratis tracks information about file systems created using Stratis that XFS is not aware of, and changes made using XFS do not automatically create updates in Stratis. Users must not reformat or reconfigure XFS file systems that are managed by Stratis.

Stratis creates links to file systems at the **/dev/stratis/my-pool/my-fs** path.

**NOTE**

Stratis uses many Device Mapper devices, which show up in **dmsetup** listings and the **/proc/partitions** file. Similarly, the **lsblk** command output reflects the internal workings and layers of Stratis.

33.2. ADDING BLOCK DEVICES TO A STRATIS POOL

This procedure adds one or more block devices to a Stratis pool to be usable by Stratis file systems.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- The block devices that you are adding to the Stratis pool are not in use and not mounted.
- The block devices that you are adding to the Stratis pool are at least 1 GiB in size each.

Procedure

- To add one or more block devices to the pool, use:

```
# stratis pool add-data my-pool device-1 device-2 device-n
```

Additional resources

- **stratis(8)** man page

33.3. ADDITIONAL RESOURCES

- [The Stratis Storage website](#)

CHAPTER 34. MONITORING STRATIS FILE SYSTEMS

As a Stratis user, you can view information about Stratis volumes on your system to monitor their state and free space.



IMPORTANT

Stratis is a Technology Preview feature only. Technology Preview features are not supported with Red Hat production service level agreements (SLAs) and might not be functionally complete. Red Hat does not recommend using them in production. These features provide early access to upcoming product features, enabling customers to test functionality and provide feedback during the development process. For more information about the support scope of Red Hat Technology Preview features, see <https://access.redhat.com/support/offerings/techpreview>.

34.1. STRATIS SIZES REPORTED BY DIFFERENT UTILITIES

This section explains the difference between Stratis sizes reported by standard utilities such as **df** and the **stratis** utility.

Standard Linux utilities such as **df** report the size of the XFS file system layer on Stratis, which is 1 TiB. This is not useful information, because the actual storage usage of Stratis is less due to thin provisioning, and also because Stratis automatically grows the file system when the XFS layer is close to full.



IMPORTANT

Regularly monitor the amount of data written to your Stratis file systems, which is reported as the *Total Physical Used* value. Make sure it does not exceed the *Total Physical Size* value.

Additional resources

- **stratis(8)** man page.

34.2. DISPLAYING INFORMATION ABOUT STRATIS VOLUMES

This procedure lists statistics about your Stratis volumes, such as the total, used, and free size or file systems and block devices belonging to a pool.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.

Procedure

- To display information about all **block devices** used for Stratis on your system:

```
# stratis blockdev
```

Pool Name	Device Node	Physical Size	State	Tier
my-pool	/dev/sdb	9.10 TiB	In-use	Data

- To display information about all Stratis **pools** on your system:

```
# stratis pool
```

Name	Total Physical Size	Total Physical Used
<i>my-pool</i>	<i>9.10 TiB</i>	<i>598 MiB</i>

- To display information about all Stratis **file systems** on your system:

```
# stratis filesystem
```

Pool Name	Name	Used	Created	Device
<i>my-pool</i>	<i>my-fs</i>	<i>546 MiB</i>	<i>Nov 08 2018 08:03</i>	<i>/dev/stratis/my-pool/my-fs</i>

Additional resources

- **stratis(8)** man page.

34.3. ADDITIONAL RESOURCES

- [The Stratis Storage website](#)

CHAPTER 35. USING SNAPSHOTS ON STRATIS FILE SYSTEMS

You can use snapshots on Stratis file systems to capture file system state at arbitrary times and restore it in the future.



IMPORTANT

Stratis is a Technology Preview feature only. Technology Preview features are not supported with Red Hat production service level agreements (SLAs) and might not be functionally complete. Red Hat does not recommend using them in production. These features provide early access to upcoming product features, enabling customers to test functionality and provide feedback during the development process. For more information about the support scope of Red Hat Technology Preview features, see <https://access.redhat.com/support/offerings/techpreview>.

35.1. CHARACTERISTICS OF STRATIS SNAPSHOTS

In Stratis, a snapshot is a regular Stratis file system created as a copy of another Stratis file system. The snapshot initially contains the same file content as the original file system, but can change as the snapshot is modified. Whatever changes you make to the snapshot will not be reflected in the original file system.

The current snapshot implementation in Stratis is characterized by the following:

- A snapshot of a file system is another file system.
- A snapshot and its origin are not linked in lifetime. A snapshotted file system can live longer than the file system it was created from.
- A file system does not have to be mounted to create a snapshot from it.
- Each snapshot uses around half a gigabyte of actual backing storage, which is needed for the XFS log.

35.2. CREATING A STRATIS SNAPSHOT

This procedure creates a Stratis file system as a snapshot of an existing Stratis file system.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis file system. See [Creating a Stratis filesystem](#).

Procedure

- To create a Stratis snapshot, use:

```
# stratis fs snapshot my-pool my-fs my-fs-snapshot
```

Additional resources

- **stratis(8)** man page.

35.3. ACCESSING THE CONTENT OF A STRATIS SNAPSHOT

This procedure mounts a snapshot of a Stratis file system to make it accessible for read and write operations.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis snapshot. See [Creating a Stratis filesystem](#).

Procedure

- To access the snapshot, mount it as a regular file system from the **/dev/stratis/my-pool/** directory:

```
# mount /dev/stratis/my-pool/my-fs-snapshot mount-point
```

Additional resources

- [Mounting a Stratis file system](#).
- **mount(8)** man page.

35.4. REVERTING A STRATIS FILE SYSTEM TO A PREVIOUS SNAPSHOT

This procedure reverts the content of a Stratis file system to the state captured in a Stratis snapshot.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis snapshot. See [Creating a Stratis snapshot](#).

Procedure

1. Optionally, back up the current state of the file system to be able to access it later:

```
# stratis filesystem snapshot my-pool my-fs my-fs-backup
```

2. Unmount and remove the original file system:

```
# umount /dev/stratis/my-pool/my-fs
# stratis filesystem destroy my-pool my-fs
```

3. Create a copy of the snapshot under the name of the original file system:

```
# stratis filesystem snapshot my-pool my-fs-snapshot my-fs
```

4. Mount the snapshot, which is now accessible with the same name as the original file system:

```
# mount /dev/stratis/my-pool/my-fs mount-point
```

The content of the file system named *my-fs* is now identical to the snapshot *my-fs-snapshot*.

Additional resources

- **stratis(8)** man page.

35.5. REMOVING A STRATIS SNAPSHOT

This procedure removes a Stratis snapshot from a pool. Data on the snapshot are lost.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis snapshot. See [Creating a Stratis snapshot](#).

Procedure

1. Unmount the snapshot:

```
# umount /dev/stratis/my-pool/my-fs-snapshot
```

2. Destroy the snapshot:

```
# stratis filesystem destroy my-pool my-fs-snapshot
```

Additional resources

- **stratis(8)** man page.

35.6. ADDITIONAL RESOURCES

- [The Stratis Storage website](#)

CHAPTER 36. REMOVING STRATIS FILE SYSTEMS

You can remove an existing Stratis file system, or a Stratis pool, by destroying data on them.



IMPORTANT

Stratis is a Technology Preview feature only. Technology Preview features are not supported with Red Hat production service level agreements (SLAs) and might not be functionally complete. Red Hat does not recommend using them in production. These features provide early access to upcoming product features, enabling customers to test functionality and provide feedback during the development process. For more information about the support scope of Red Hat Technology Preview features, see <https://access.redhat.com/support/offerings/techpreview>.

36.1. COMPONENTS OF A STRATIS VOLUME

Learn about the components that comprise a Stratis volume.

Externally, Stratis presents the following volume components in the command-line interface and the API:

blockdev

Block devices, such as a disk or a disk partition.

pool

Composed of one or more block devices.

A pool has a fixed total size, equal to the size of the block devices.

The pool contains most Stratis layers, such as the non-volatile data cache using the **dm-cache** target.

Stratis creates a **/dev/stratis/my-pool/** directory for each pool. This directory contains links to devices that represent Stratis file systems in the pool.

filesystem

Each pool can contain one or more file systems, which store files.

File systems are thinly provisioned and do not have a fixed total size. The actual size of a file system grows with the data stored on it. If the size of the data approaches the virtual size of the file system, Stratis grows the thin volume and the file system automatically.

The file systems are formatted with XFS.



IMPORTANT

Stratis tracks information about file systems created using Stratis that XFS is not aware of, and changes made using XFS do not automatically create updates in Stratis. Users must not reformat or reconfigure XFS file systems that are managed by Stratis.

Stratis creates links to file systems at the **/dev/stratis/my-pool/my-fs** path.

**NOTE**

Stratis uses many Device Mapper devices, which show up in **dmsetup** listings and the **/proc/partitions** file. Similarly, the **lsblk** command output reflects the internal workings and layers of Stratis.

36.2. REMOVING A STRATIS FILE SYSTEM

This procedure removes an existing Stratis file system. Data stored on it are lost.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis file system. See [Creating a Stratis filesystem](#).

Procedure

1. Unmount the file system:

```
# umount /dev/stratis/my-pool/my-fs
```

2. Destroy the file system:

```
# stratis filesystem destroy my-pool my-fs
```

3. Verify that the file system no longer exists:

```
# stratis filesystem list my-pool
```

Additional resources

- **stratis(8)** man page.

36.3. REMOVING A STRATIS POOL

This procedure removes an existing Stratis pool. Data stored on it are lost.

Prerequisites

- Stratis is installed. See [Installing Stratis](#).
- The **stratisd** service is running.
- You have created a Stratis pool:
 - To create an unencrypted pool, see [Creating an unencrypted Stratis pool](#)
 - To create an encrypted pool, see [Creating an encrypted Stratis pool](#).

Procedure

1. List file systems on the pool:

```
# stratis filesystem list my-pool
```

2. Unmount all file systems on the pool:

```
# umount /dev/stratis/my-pool/my-fs-1 \  
         /dev/stratis/my-pool/my-fs-2 \  
         /dev/stratis/my-pool/my-fs-n
```

3. Destroy the file systems:

```
# stratis filesystem destroy my-pool my-fs-1 my-fs-2
```

4. Destroy the pool:

```
# stratis pool destroy my-pool
```

5. Verify that the pool no longer exists:

```
# stratis pool list
```

Additional resources

- **stratis(8)** man page.

36.4. ADDITIONAL RESOURCES

- [The Stratis Storage website](#)

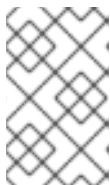
CHAPTER 37. GETTING STARTED WITH AN EXT4 FILE SYSTEM

As a system administrator, you can create, mount, resize, backup, and restore an ext4 file system. The ext4 file system is a scalable extension of the ext3 file system. With Red Hat Enterprise Linux 9, it can support a maximum individual file size of **16** terabytes, and file system to a maximum of **50** terabytes.

37.1. FEATURES OF AN EXT4 FILE SYSTEM

Following are the features of an ext4 file system:

- Using extents: The ext4 file system uses extents, which improves performance when using large files and reduces metadata overhead for large files.
- Ext4 labels unallocated block groups and inode table sections accordingly, which allows the block groups and table sections to be skipped during a file system check. It leads to a quick file system check, which becomes more beneficial as the file system grows in size.
- Metadata checksum: By default, this feature is enabled in Red Hat Enterprise Linux 9.
- Allocation features of an ext4 file system:
 - Persistent pre-allocation
 - Delayed allocation
 - Multi-block allocation
 - Stripe-aware allocation
- Extended attributes (**xattr**): This allows the system to associate several additional name and value pairs per file.
- Quota journaling: This avoids the need for lengthy quota consistency checks after a crash.



NOTE

The only supported journaling mode in ext4 is **data=ordered** (default). For more information, see [Is the EXT journaling option "data=writeback" supported in RHEL? Knowledgebase article](#).

- Subsecond timestamps – This gives timestamps to the subsecond.

Additional resources

- **ext4** man page.

37.2. CREATING AN EXT4 FILE SYSTEM

As a system administrator, you can create an ext4 file system on a block device using **mkfs.ext4** command.

Prerequisites

- A partition on your disk. For information about creating MBR or GPT partitions, see [Creating a partition table on a disk with parted](#).
- Alternatively, use an LVM or MD volume.

Procedure

1. To create an ext4 file system:

- For a regular-partition device, an LVM volume, an MD volume, or a similar device, use the following command:

```
# mkfs.ext4 /dev/block_device
```

Replace `/dev/block_device` with the path to a block device.

For example, `/dev/sdb1`, `/dev/disk/by-uuid/05e99ec8-def1-4a5e-8a9d-5945339ceb2a`, or `/dev/my-volgroup/my-lv`. In general, the default options are optimal for most usage scenarios.

- For striped block devices (for example, RAID5 arrays), the stripe geometry can be specified at the time of file system creation. Using proper stripe geometry enhances the performance of an ext4 file system. For example, to create a file system with a 64k stride (that is, 16 x 4096) on a 4k-block file system, use the following command:

```
# mkfs.ext4 -E stride=16,stripe-width=64 /dev/block_device
```

In the given example:

- `stride=value`: Specifies the RAID chunk size
- `stripe-width=value`: Specifies the number of data disks in a RAID device, or the number of stripe units in the stripe.



NOTE

- To specify a UUID when creating a file system:

```
# mkfs.ext4 -U UUID /dev/block_device
```

Replace `UUID` with the UUID you want to set: for example, **7cd65de3-e0be-41d9-b66d-96d749c02da7**.

Replace `/dev/block_device` with the path to an ext4 file system to have the UUID added to it: for example, `/dev/sda8`.

- To specify a label when creating a file system:

```
# mkfs.ext4 -L label-name /dev/block_device
```

2. To view the created ext4 file system:

```
# blkid
```

Additional resources

- **ext4** man page.
- **mkfs.ext4** man page.

37.3. MOUNTING AN EXT4 FILE SYSTEM

As a system administrator, you can mount an ext4 file system using the **mount** utility.

Prerequisites

- An ext4 file system. For information about creating an ext4 file system, see [Creating an ext4 file system](#).

Procedure

1. To create a mount point to mount the file system:

```
# mkdir /mount/point
```

Replace */mount/point* with the directory name where mount point of the partition must be created.

2. To mount an ext4 file system:

- To mount an ext4 file system with no extra options:

```
# mount /dev/block_device /mount/point
```

- To mount the file system persistently, see [Persistently mounting file systems](#).

3. To view the mounted file system:

```
# df -h
```

Additional resources

- **mount** man page.
- **ext4** man page.
- **fstab** man page.
- [Mounting file systems](#).

37.4. RESIZING AN EXT4 FILE SYSTEM

As a system administrator, you can resize an ext4 file system using the **resize2fs** utility. The **resize2fs** utility reads the size in units of file system block size, unless a suffix indicating a specific unit is used. The following suffixes indicate specific units:

- s (sectors) - **512** byte sectors

- K (kilobytes) - **1,024** bytes
- M (megabytes) - **1,048,576** bytes
- G (gigabytes) - **1,073,741,824** bytes
- T (terabytes) - **1,099,511,627,776** bytes

Prerequisites

- An ext4 file system. For information about creating an ext4 file system, see [Creating an ext4 file system](#).
- An underlying block device of an appropriate size to hold the file system after resizing.

Procedure

1. To resize an ext4 file system, take the following steps:

- To shrink and grow the size of an unmounted ext4 file system:

```
# umount /dev/block_device
# e2fsck -f /dev/block_device
# resize2fs /dev/block_device size
```

Replace `/dev/block_device` with the path to the block device, for example `/dev/sdb1`.

Replace `size` with the required resize value using **s**, **K**, **M**, **G**, and **T** suffixes.

- An ext4 file system may be grown while mounted using the **resize2fs** command:

```
# resize2fs /mount/device size
```



NOTE

The size parameter is optional (and often redundant) when expanding. The **resize2fs** automatically expands to fill the available space of the container, usually a logical volume or partition.

2. To view the resized file system:

```
# df -h
```

Additional resources

- **resize2fs** man page.
- **e2fsck** man page.
- **ext4** man page.

37.5. COMPARISON OF TOOLS USED WITH EXT4 AND XFS

This section compares which tools to use to accomplish common tasks on the ext4 and XFS file systems.

Task	ext4	XFS
Create a file system	mkfs.ext4	mkfs.xfs
File system check	e2fsck	xfs_repair
Resize a file system	resize2fs	xfs_growfs
Save an image of a file system	e2image	xfs_metadump and xfs_mdrestore
Label or tune a file system	tune2fs	xfs_admin
Back up a file system	tar and rsync	xfsdump and xfsrestore
Quota management	quota	xfs_quota
File mapping	filefrag	xfs_bmap



NOTE

If you want a complete client-server solution for backups over network, you can use **bacula** backup utility that is available in RHEL 9. For more information about Bacula, see [Bacula backup solution](#).