

Veritas™ File System Administrator's Guide

HP-UX

5.0.1

Veritas File System Administrator's Guide

The software described in this book is furnished under a license agreement and may be used only in accordance with the terms of the agreement.

Product version: 5.0.1

Documentation version: 5.0.1.0

Legal Notice

Copyright © 2009 Symantec Corporation. All rights reserved.

Symantec, the Symantec Logo, Veritas, and Veritas Storage Foundation are trademarks or registered trademarks of Symantec Corporation or its affiliates in the U.S. and other countries. Other names may be trademarks of their respective owners.

The product described in this document is distributed under licenses restricting its use, copying, distribution, and decompilation/reverse engineering. No part of this document may be reproduced in any form by any means without prior written authorization of Symantec Corporation and its licensors, if any.

THE DOCUMENTATION IS PROVIDED "AS IS" AND ALL EXPRESS OR IMPLIED CONDITIONS, REPRESENTATIONS AND WARRANTIES, INCLUDING ANY IMPLIED WARRANTY OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE OR NON-INFRINGEMENT, ARE DISCLAIMED, EXCEPT TO THE EXTENT THAT SUCH DISCLAIMERS ARE HELD TO BE LEGALLY INVALID. SYMANTEC CORPORATION SHALL NOT BE LIABLE FOR INCIDENTAL OR CONSEQUENTIAL DAMAGES IN CONNECTION WITH THE FURNISHING, PERFORMANCE, OR USE OF THIS DOCUMENTATION. THE INFORMATION CONTAINED IN THIS DOCUMENTATION IS SUBJECT TO CHANGE WITHOUT NOTICE.

The Licensed Software and Documentation are deemed to be commercial computer software as defined in FAR 12.212 and subject to restricted rights as defined in FAR Section 52.227-19 "Commercial Computer Software - Restricted Rights" and DFARS 227.7202, "Rights in Commercial Computer Software or Commercial Computer Software Documentation", as applicable, and any successor regulations. Any use, modification, reproduction release, performance, display or disclosure of the Licensed Software and Documentation by the U.S. Government shall be solely in accordance with the terms of this Agreement.

Symantec Corporation
350 Ellis Street
Mountain View, CA 94043
<http://www.symantec.com>

Technical Support

Symantec Technical Support maintains support centers globally. Technical Support's primary role is to respond to specific queries about product features and functionality. The Technical Support group also creates content for our online Knowledge Base. The Technical Support group works collaboratively with the other functional areas within Symantec to answer your questions in a timely fashion. For example, the Technical Support group works with Product Engineering and Symantec Security Response to provide alerting services and virus definition updates.

Symantec's maintenance offerings include the following:

- A range of support options that give you the flexibility to select the right amount of service for any size organization
- Telephone and Web-based support that provides rapid response and up-to-the-minute information
- Upgrade assurance that delivers automatic software upgrade protection
- Global support that is available 24 hours a day, 7 days a week
- Advanced features, including Account Management Services

For information about Symantec's Maintenance Programs, you can visit our Web site at the following URL:

www.symantec.com/techsupp/

Contacting Technical Support

Customers with a current maintenance agreement may access Technical Support information at the following URL:

www.symantec.com/techsupp/

Before contacting Technical Support, make sure you have satisfied the system requirements that are listed in your product documentation. Also, you should be at the computer on which the problem occurred, in case it is necessary to replicate the problem.

When you contact Technical Support, please have the following information available:

- Product release level
- Hardware information
- Available memory, disk space, and NIC information
- Operating system

- Version and patch level
- Network topology
- Router, gateway, and IP address information
- Problem description:
 - Error messages and log files
 - Troubleshooting that was performed before contacting Symantec
 - Recent software configuration changes and network changes

Licensing and registration

If your Symantec product requires registration or a license key, access our technical support Web page at the following URL:

www.symantec.com/techsupp/

Customer service

Customer service information is available at the following URL:

www.symantec.com/techsupp/

Customer Service is available to assist with the following types of issues:

- Questions regarding product licensing or serialization
- Product registration updates, such as address or name changes
- General product information (features, language availability, local dealers)
- Latest information about product updates and upgrades
- Information about upgrade assurance and maintenance contracts
- Information about the Symantec Buying Programs
- Advice about Symantec's technical support options
- Nontechnical presales questions
- Issues that are related to CD-ROMs or manuals

Maintenance agreement resources

If you want to contact Symantec regarding an existing maintenance agreement, please contact the maintenance agreement administration team for your region as follows:

Asia-Pacific and Japan	customercare_apac@symantec.com
Europe, Middle-East, and Africa	semea@symantec.com
North America and Latin America	supportsolutions@symantec.com

Additional enterprise services

Symantec offers a comprehensive set of services that allow you to maximize your investment in Symantec products and to develop your knowledge, expertise, and global insight, which enable you to manage your business risks proactively.

Enterprise services that are available include the following:

Symantec Early Warning Solutions	These solutions provide early warning of cyber attacks, comprehensive threat analysis, and countermeasures to prevent attacks before they occur.
Managed Security Services	These services remove the burden of managing and monitoring security devices and events, ensuring rapid response to real threats.
Consulting Services	Symantec Consulting Services provide on-site technical expertise from Symantec and its trusted partners. Symantec Consulting Services offer a variety of prepackaged and customizable options that include assessment, design, implementation, monitoring, and management capabilities. Each is focused on establishing and maintaining the integrity and availability of your IT resources.
Educational Services	Educational Services provide a full array of technical training, security education, security certification, and awareness communication programs.

To access more information about Enterprise services, please visit our Web site at the following URL:

www.symantec.com

Select your country or language from the site index.

Contents

Technical Support	4
Chapter 1 Introducing Veritas File System	15
About Veritas File System	15
Logging	16
Extents	16
File system disk layouts	16
Veritas File System features	16
Extent-based allocation	18
Extent attributes	20
Fast file system recovery	20
Extended mount options	21
Enhanced data integrity modes	22
Enhanced performance mode	23
Temporary file system mode	23
Improved synchronous writes	23
Support for large files	23
Access Control Lists	24
Storage Checkpoints	24
Online backup	24
Quotas	25
Support for databases	25
Cluster file systems	25
Cross-platform data sharing	26
File Change Log	26
Multi-volume support	26
Dynamic Storage Tiering	27
Thin Reclamation of a file system	27
Veritas File System performance enhancements	27
About enhanced I/O performance	28
Using Veritas File System	29
Online system administration	29
Application program interface	30

Chapter 2	VxFS performance: creating, mounting, and tuning file systems	33
	Creating a VxFS file system	33
	Block size	33
	Intent log size	34
	Mounting a VxFS file system	34
	The log mode	35
	The delaylog mode	36
	The tmplog mode	36
	The logiosize mode	37
	The nodatainlog mode	37
	The blkclear mode	38
	The mincache mode	38
	The convosync mode	39
	The ioerror mode	40
	The largefiles nolargefiles option	42
	The cio option	43
	The mntlock mntunlock option	43
	Combining mount command options	44
	Tuning the VxFS file system	44
	Tuning inode table size	45
	VxFS buffer cache high water mark	46
	Number of links to a file	47
	VxFS inode free time lag	48
	Monitoring free space	48
	Monitoring fragmentation	49
	Thin Reclamation	50
	Tuning I/O	51
	Tuning VxFS I/O parameters	51
	Tunable I/O parameters	52
	File system tuning guidelines	60
Chapter 3	Extent attributes	63
	About extent attributes	63
	Reservation: preallocating space to a file	64
	Fixed extent size	64
	Other controls	65
	Commands related to extent attributes	66
	Failure to preserve extent attributes	67

Chapter 4	VxFS I/O Overview	69
	About VxFS I/O	69
	Buffered and Direct I/O	70
	Direct I/O	70
	Unbuffered I/O	71
	Data synchronous I/O	71
	Concurrent I/O	72
	Cache advisories	73
	Freezing and thawing a file system	73
	Getting the I/O size	74
Chapter 5	Online backup using file system snapshots	75
	About snapshot file systems	75
	Snapshot file system backups	76
	Creating a snapshot file system	77
	Backup examples	77
	Snapshot file system performance	78
	Differences Between Snapshots and Storage Checkpoints	79
	About snapshot file system disk structure	79
	How a snapshot file system works	80
Chapter 6	Storage Checkpoints	83
	About Storage Checkpoints	83
	How Storage Checkpoints differ from snapshots	84
	How a Storage Checkpoint works	85
	Copy-on-write	87
	Types of Storage Checkpoints	88
	Data Storage Checkpoints	88
	nodata Storage Checkpoints	88
	Removable Storage Checkpoints	89
	Non-mountable Storage Checkpoints	89
	Storage Checkpoint administration	90
	Creating a Storage Checkpoint	91
	Removing a Storage Checkpoint	92
	Accessing a Storage Checkpoint	92
	Converting a data Storage Checkpoint to a nodata Storage Checkpoint	94
	Space management considerations	101
	Restoring a file system from a Storage Checkpoint	102
	Restoring a file from a Storage Checkpoint	102
	Storage Checkpoint quotas	107

Chapter 7	Quotas	109
	About quota limits	109
	About quota files on Veritas File System	110
	About quota commands	110
	Using quotas	110
	Turning on quotas	111
	Turning on quotas at mount time	111
	Editing user quotas	111
	Modifying time limits	111
	Viewing disk quotas and usage	112
	Displaying blocks owned by users or groups	112
	Turning off quotas	112
Chapter 8	File Change Log	113
	About File Change Log	113
	About the File Change Log file	114
	File Change Log administrative interface	115
	File Change Log programmatic interface	117
	Summary of API functions	119
	Reverse path name lookup	120
Chapter 9	Multi-volume file systems	121
	About multi-volume support	122
	About volume types	122
	Features implemented using multi-volume support	122
	Volume availability	123
	About volume sets	124
	Creating and managing volume sets	124
	Creating multi-volume file systems	125
	Example of creating a multi-volume file system	126
	Converting a single volume file system to a multi-volume file system	127
	Removing a volume from a multi-volume file system	128
	Forcibly removing a volume	129
	Moving volume 0	129
	About allocation policies	129
	Assigning allocation policies	129
	Querying allocation policies	131
	Assigning pattern tables to directories	132
	Assigning pattern tables to file systems	132
	Allocating data	133

Volume encapsulation	134
Encapsulating a volume	134
Deencapsulating a volume	136
Reporting file extents	136
Examples of reporting file extents	137
Load balancing	138
Defining and assigning a load balancing allocation policy	138
Rebalancing extents	138
Converting a multi-volume file system to a single volume file system	139
Converting to a single volume file system	139
 Chapter 10 Dynamic Storage Tiering	 143
About Dynamic Storage Tiering	143
Placement classes	145
Tagging volumes as placement classes	146
Listing placement classes	146
Administering placement policies	146
Assigning a placement policy	147
Unassigning a placement policy	147
Analyzing the space impact of enforcing a placement policy	147
Querying which files will be affected by enforcing a placement policy	148
Enforcing a placement policy	148
Validating a placement policy	149
File placement policy grammar	150
File placement policy rules	150
SELECT statement	151
CREATE statement	154
RELOCATE statement	156
DELETE statement	167
Calculating I/O temperature and access temperature	169
Multiple criteria in file placement policy rule statements	173
Multiple file selection criteria in SELECT statement clauses	173
Multiple placement classes in <ON> clauses of CREATE statements and in <TO> clauses of RELOCATE statements	174
Multiple placement classes in <FROM> clauses of RELOCATE and DELETE statements	175
Multiple conditions in <WHEN> clauses of RELOCATE and DELETE statements	176

File placement policy rule and statement ordering	176
File placement policies and extending files	179

Appendix A	Quick Reference	181
	Command summary	181
	Online manual pages	184
	Creating a VxFS file system	190
	Example of creating a file system	191
	Converting a file system to VxFS	192
	Example of converting a file system	192
	Mounting a file system	192
	Mount options	193
	Example of mounting a file system	194
	Editing the fstab file	195
	Unmounting a file system	195
	Example of unmounting a file system	196
	Displaying information on mounted file systems	196
	Example of displaying information on mounted file systems	196
	Identifying file system types	197
	Example of determining a file system's type	198
	Resizing a file system	198
	Extending a file system using fsadm	198
	Shrinking a file system	199
	Reorganizing a file system	200
	Extending a file system using extendfs	201
	Backing up and restoring a file system	202
	Creating and mounting a snapshot file system	202
	Backing up a file system	203
	Restoring a file system	204
	Using quotas	204
	Turning on quotas	205
	Setting up user quotas	205
	Viewing quotas	206
	Turning off quotas	206

Appendix B	Diagnostic messages	209
	File system response to problems	209
	Recovering a disabled file system	210
	About kernel messages	210
	About global message IDs	210
	Kernel messages	211

	About unique message identifiers	250
	Unique message identifiers	251
Appendix C	Disk layout	255
	About disk layouts	255
	Supported disk layouts and operating systems	256
	VxFS Version 4 disk layout	258
	VxFS Version 5 disk layout	261
	VxFS Version 6 disk layout	262
	VxFS Version 7 disk layout	263
Glossary		265
Index		271

Introducing Veritas File System

This chapter includes the following topics:

- [About Veritas File System](#)
- [Veritas File System features](#)
- [Veritas File System performance enhancements](#)
- [Using Veritas File System](#)

About Veritas File System

A file system is simply a method for storing and organizing computer files and the data they contain to make it easy to find and access them. More formally, a file system is a set of abstract data types (such as metadata) that are implemented for the storage, hierarchical organization, manipulation, navigation, access, and retrieval of data.

Veritas File System (VxFS) was the first commercial journaling file system. With journaling, metadata changes are first written to a log (or journal) then to disk. Since changes do not need to be written in multiple places, throughput is much faster as the metadata is written asynchronously.

VxFS is also an extent-based, intent logging file system. VxFS is designed for use in operating environments that require high performance and availability and deal with large amounts of data.

VxFS major components include:

- Logging
- Extents

- File system disk layouts

Logging

A key aspect of any file system is how to recover if a system crash occurs. Earlier methods required a time-consuming scan of the entire file system. A better solution is the method of logging (or journaling) the metadata of files.

VxFS logs new attribute information into a reserved area of the file system, whenever file system changes occur. The file system writes the actual data to disk only after the write of the metadata to the log is complete. If and when a system crash occurs, the system recovery code analyzes the metadata log and tries to clean up only those files. Without logging, a file system check (`fsck`) must look at all of the metadata.

Intent logging minimizes system downtime after abnormal shutdowns by logging file system transactions. When the system is halted unexpectedly, this log can be replayed and outstanding transactions completed. The check and repair time for file systems can be reduced to a few seconds, regardless of the file system size.

By default, VxFS file systems log file transactions before they are committed to disk, reducing time spent checking and repairing file systems after the system is halted unexpectedly.

Extents

An extent is a contiguous area of storage in a computer file system, reserved for a file. When starting to write to a file, a whole extent is allocated. When writing to the file again, the data continues where the previous write left off. This reduces or eliminates file fragmentation.

Since VxFS is an extent-based file system, addressing is done through extents (which can consist of multiple blocks) rather than in single-block segments. Extents can therefore enhance file system throughput.

File system disk layouts

The disk layout is the way file system information is stored on disk. On VxFS, several disk layout versions, numbered 1 through 7, were created to support various new features and specific UNIX environments. Currently, only the Version 4, 5, 6, and 7 disk layouts can be created and mounted.

Veritas File System features

VxFS includes the following features:

- **Extent-based allocation**
Extents allow disk I/O to take place in units of multiple blocks if storage is allocated in consecutive blocks.
- **Extent attributes**
Extent attributes are the extent allocation policies associated with a file.
- **Fast file system recovery**
VxFS provides fast recovery of a file system from system failure.
- **Extended mount options**
The VxFS file system supports extended `mount` options to specify enhanced data integrity modes, enhanced performance modes, temporary file system modes, improved synchronous writes, and large file sizes.
- **Enhanced performance mode**
VxFS provides mount options to improve performance.
- **Large files and file systems support**
VxFS supports files larger than two gigabytes and large file systems up to 256 terabytes.
- **Online backup**
VxFS provides online data backup using the snapshot feature.
- **Quotas**
VxFS supports quotas, which allocate per-user and per-group quotas and limit the use of two principal resources: files and data blocks.
- **Cluster File System**
Clustered file systems are an extension of VxFS that support concurrent direct media access from multiple systems.
- **Improved database performance**
- **Storage Checkpoints**
Backup and restore applications can leverage Storage Checkpoint, a disk- and I/O-efficient copying technology for creating periodic frozen images of a file system.
- **Cross-platform data sharing**
Cross-platform data sharing allows data to be serially shared among heterogeneous systems where each system has direct access to the physical devices that hold the data.
- **File Change Log**
The VxFS File Change Log tracks changes to files and directories in a file system.

- **Multi-volume support**
The multi-volume support feature allows several volumes to be represented by a single logical object.
- **Dynamic Storage Tiering**
The Dynamic Storage Tiering (DST) option allows you to configure policies that automatically relocate files from one volume to another, or relocate files by running file relocation commands, which can improve performance for applications that access specific types of files.
- **Storage Foundation Thin Reclamation**
The Thin Reclamation feature allows you to release free data blocks of a VxFS file system to the free storage pool of a Thin Storage LUN. This feature is only supported on file systems mounted on a VxVM volume.

Note: VxFS supports all HFS file system features and facilities except for the linking, removing, or renaming of “.” and “.” directory entries. Such operations may disrupt file system sanity.

Extent-based allocation

Disk space is allocated in 1024-byte sectors to form logical blocks. VxFS supports logical block sizes of 1024, 2048, 4096, and 8192 bytes. The default block size is the larger of 1K or the device's hardware sector size for file system sizes of up to 1 TB, and 8K for file system sizes 1 TB or larger.

An extent is defined as one or more adjacent blocks of data within the file system. An extent is presented as an address-length pair, which identifies the starting block address and the length of the extent (in file system or logical blocks). VxFS allocates storage in groups of extents rather than a block at a time.

Extents allow disk I/O to take place in units of multiple blocks if storage is allocated in consecutive blocks. For sequential I/O, multiple block operations are considerably faster than block-at-a-time operations; almost all disk drives accept I/O operations of multiple blocks.

Extent allocation only slightly alters the interpretation of addressed blocks from the inode structure compared to block based inodes. A VxFS inode references 10 direct extents, each of which are pairs of starting block addresses and lengths in blocks.

The VxFS inode supports different types of extents, namely `ext4` and `typed`. Inodes with `ext4` extents also point to two indirect address extents, which contain the addresses of first and second extents:

first	Used for single indirection. Each entry in the extent indicates the starting block number of an indirect data extent
second	Used for double indirection. Each entry in the extent indicates the starting block number of a single indirect address extent.

Each indirect address extent is 8K long and contains 2048 entries. All indirect data extents for a file must be the same size; this size is set when the first indirect data extent is allocated and stored in the inode. Directory inodes always use an 8K indirect data extent size. By default, regular file inodes also use an 8K indirect data extent size that can be altered with `vxtunefs`; these inodes allocate the indirect data extents in clusters to simulate larger extents.

Typed extents

VxFS has an inode block map organization for indirect extents known as `typed` extents. Each entry in the block map has a typed descriptor record containing a type, offset, starting block, and number of blocks.

Indirect and data extents use this format to identify logical file offsets and physical disk locations of any given extent.

The extent descriptor fields are defined as follows:

type	Identifies uniquely an extent descriptor record and defines the record's length and format.
offset	Represents the logical file offset in blocks for a given descriptor. Used to optimize lookups and eliminate hole descriptor entries.
starting block	Is the starting file system block of the extent.
number of blocks	Is the number of contiguous blocks in the extent.

`Typed` extents have the following characteristics:

- Indirect address blocks are fully typed and may have variable lengths up to a maximum and optimum size of 8K. On a fragmented file system, indirect extents may be smaller than 8K depending on space availability. VxFS always tries to obtain 8K indirect extents but resorts to smaller indirects if necessary.
- Indirect data extents are variable in size to allow files to allocate large, contiguous extents and take full advantage of optimized I/O in VxFS.
- Holes in sparse files require no storage and are eliminated by typed records. A hole is determined by adding the offset and length of a descriptor and comparing the result with the offset of the next record.

- While there are no limits on the levels of indirection, lower levels are expected in this format since data extents have variable lengths.
- This format uses a type indicator that determines its record format and content and accommodates new requirements and functionality for future types.

The current typed format is used on regular files and directories only when indirection is needed. Typed records are longer than the previous format and require less direct entries in the inode. Newly created files start out using the old format, which allows for ten direct extents in the inode. The inode's block map is converted to the typed format when indirection is needed to offer the advantages of both formats.

Extent attributes

VxFS allocates disk space to files in groups of one or more extents. VxFS also allows applications to control some aspects of the extent allocation. Extent attributes are the extent allocation policies associated with a file.

The `setext` and `getext` commands allow the administrator to set or view extent attributes associated with a file, as well as to preallocate space for a file.

See the `setext(1M)` and `getext(1M)` manual pages.

The `vxtunefs` command allows the administrator to set or view the default indirect data extent size of a file system.

See the `vxtunefs(1M)` manual page.

Fast file system recovery

Most file systems rely on full structural verification by the `fsck` utility as the only means to recover from a system failure. For large disk configurations, this involves a time-consuming process of checking the entire structure, verifying that the file system is intact, and correcting any inconsistencies. VxFS provides fast recovery with the VxFS intent log and VxFS intent log resizing features.

VxFS intent log

VxFS reduces system failure recovery times by tracking file system activity in the VxFS intent log. This feature records pending changes to the file system structure in a circular intent log. The intent log recovery feature is not readily apparent to users or a system administrator except during a system failure. During system failure recovery, the VxFS `fsck` utility performs an intent log replay, which scans the intent log and nullifies or completes file system operations that were active when the system failed. The file system can then be mounted without completing

a full structural check of the entire file system. Replaying the intent log may not completely recover the damaged file system structure if there was a disk hardware failure; hardware problems may require a complete system check using the `fsck` utility provided with VxFS.

See [“The log option and data integrity”](#) on page 22.

VxFS intent log resizing

The VxFS intent log is allocated when the file system is first created. The size of the intent log is based on the size of the file system—the larger the file system, the larger the intent log. The maximum default intent log size for disk layout Versions 4, 5, and 6 is 16 megabytes. The maximum default intent log size for disk layout Version 7 is 64 megabytes.

With the Version 6 and 7 disk layouts, you can dynamically increase or decrease the intent log size using the `logsize` option of the `fsadm` command. Increasing the size of the intent log can improve system performance because it reduces the number of times the log wraps around. However, increasing the intent log size can lead to greater times required for a log replay if there is a system failure.

Note: Inappropriate sizing of the intent log can have a negative impact on system performance.

See the `mkfs_vxfs(1M)` and the `fsadm_vxfs(1M)` manual pages.

Extended mount options

The VxFS file system provides the following enhancements to the `mount` command:

- Enhanced data integrity modes
- Enhanced performance mode
- Temporary file system mode
- Improved synchronous writes
- Support for large file sizes

See [“Mounting a VxFS file system”](#) on page 34.

See the `mount_vxfs(1M)` manual page.

Enhanced data integrity modes

For most UNIX file systems, including VxFS, the default mode for writing to a file is delayed, or buffered, meaning that the data to be written is copied to the file system cache and later flushed to disk.

A delayed write provides much better performance than synchronously writing the data to disk. However, in the event of a system failure, data written shortly before the failure may be lost since it was not flushed to disk. In addition, if space was allocated to the file as part of the write request, and the corresponding data was not flushed to disk before the system failure occurred, uninitialized data can appear in the file.

For the most common type of write, delayed extending writes (a delayed write that increases the file size), VxFS avoids the problem of uninitialized data appearing in the file by waiting until the data has been flushed to disk before updating the new file size to disk. If a system failure occurs before the data has been flushed to disk, the file size has not yet been updated to be uninitialized data, thus no uninitialized data appears in the file. The unused blocks that were allocated are reclaimed.

The `blkclear` option and data integrity

In environments where performance is more important than absolute data integrity, the preceding situation is not of great concern. However, VxFS supports environments that emphasize data integrity by providing the `mount -o blkclear` option that ensures uninitialized data does not appear in a file.

The `closesync` option and data integrity

VxFS provides the `mount -o mincache=closesync` option, which is useful in desktop environments with users who are likely to shut off the power on machines without halting them first. In `closesync` mode, only files that are written during the system crash or shutdown can lose data. Any changes to a file are flushed to disk when the file is closed.

The `log` option and data integrity

File systems are typically asynchronous in that structural changes to the file system are not immediately written to disk, which provides better performance. However, recent changes made to a system can be lost if a system failure occurs. Specifically, attribute changes to files and recently created files may disappear.

The `mount -o log` intent logging option guarantees that all structural changes to the file system are logged to disk before the system call returns to the application. With this option, the `rename(2)` system call flushes the source file to

disk to guarantee the persistence of the file data before renaming it. The `rename()` call is also guaranteed to be persistent when the system call returns. The changes to file system data and metadata caused by the `fsync(2)` and `fdatasync(2)` system calls are guaranteed to be persistent once the calls return.

Enhanced performance mode

VxFS has a mount option that improves performance: `delaylog`.

The `delaylog` option and enhanced performance

The default VxFS logging mode, `mount -o delaylog`, increases performance by delaying the logging of some structural changes. However, `delaylog` does not provide the equivalent data integrity as the previously described modes because recent changes may be lost during a system failure. This option provides at least the same level of data accuracy that traditional UNIX file systems provide for system failures, along with fast file system recovery.

Temporary file system mode

On most UNIX systems, temporary file system directories, such as `/tmp` and `/usr/tmp`, often hold files that do not need to be retained when the system reboots. The underlying file system does not need to maintain a high degree of structural integrity for these temporary directories. VxFS provides the `mount -o tmplog` option, which allows the user to achieve higher performance on temporary file systems by delaying the logging of most operations.

Improved synchronous writes

VxFS provides superior performance for synchronous write applications. The `mount -o datainlog` option greatly improves the performance of small synchronous writes.

The `mount -o convosync=dsync` option improves the performance of applications that require synchronous data writes but not synchronous inode time updates.

Warning: The use of the `-o convosync=dsync` option violates POSIX semantics.

Support for large files

With VxFS, you can create, mount, and manage file systems containing large files (files larger than two gigabytes).

Warning: Some applications and utilities may not work on large files.

Access Control Lists

An Access Control List (ACL) stores a series of entries that identify specific users or groups and their access privileges for a directory or file. A file may have its own ACL or may share an ACL with other files. ACLs have the advantage of specifying detailed access permissions for multiple users and groups. For VxFS file systems created with the Version 5, 6, or 7 disk layouts, up to 1024 ACL entries can be specified. ACLs are also supported on cluster file systems.

See the `getacl(1)` and `setacl(1)` manual pages.

Storage Checkpoints

To increase availability, recoverability, and performance, Veritas File System offers on-disk and online backup and restore capabilities that facilitate frequent and efficient backup strategies. Backup and restore applications can leverage a Storage Checkpoint, a disk- and I/O-efficient copying technology for creating periodic frozen images of a file system. Storage Checkpoints present a view of a file system at a point in time, and subsequently identifies and maintains copies of the original file system blocks. Instead of using a disk-based mirroring method, Storage Checkpoints save disk space and significantly reduce I/O overhead by using the free space pool available to a file system.

Storage Checkpoint functionality is separately licensed.

Online backup

VxFS provides online data backup using the snapshot feature. An image of a mounted file system instantly becomes an exact read-only copy of the file system at a specific point in time. The original file system is called the snapped file system, the copy is called the snapshot.

When changes are made to the snapped file system, the old data is copied to the snapshot. When the snapshot is read, data that has not changed is read from the snapped file system, changed data is read from the snapshot.

Backups require one of the following methods:

- Copying selected files from the snapshot file system (using `find` and `cpio`)
- Backing up the entire file system (using `fscat`)
- Initiating a full or incremental backup (using `vxdump`)

See [“About snapshot file systems”](#) on page 75.

Quotas

VxFS supports quotas, which allocate per-user quotas and limit the use of two principal resources: files and data blocks. You can assign quotas for each of these resources. Each quota consists of two limits for each resource: hard limit and soft limit.

The hard limit represents an absolute limit on data blocks or files. A user can never exceed the hard limit under any circumstances.

The soft limit is lower than the hard limit and can be exceeded for a limited amount of time. This allows users to exceed limits temporarily as long as they fall under those limits before the allotted time expires.

See [“About quota limits”](#) on page 109.

Support for databases

Databases are usually created on file systems to simplify backup, copying, and moving tasks and are slower compared to databases on raw disks.

Using Quick I/O for Databases feature with VxFS lets systems retain the benefits of having a database on a file system without sacrificing performance. Veritas Quick I/O creates regular, preallocated files to use as character devices. Databases can be created on the character devices to achieve the same performance as databases created on raw disks.

Treating regular VxFS files as raw devices has the following advantages for databases:

- Commercial database servers such as Oracle Server can issue kernel supported asynchronous I/O calls (through the `asyncdsk` or Posix AIO interface) on these pseudo devices but not on regular files.

Cluster file systems

Veritas Storage Foundation Cluster File System (SFCFS) allows clustered servers to mount and use a file system simultaneously as if all applications using the file system were running on the same server. The Veritas Volume Manager cluster functionality (CVM) makes logical volumes and raw device applications accessible through a cluster.

Beginning with SFCFS 5.0, SFCFS uses a symmetric architecture in which all nodes in the cluster can simultaneously function as metadata servers. SFCFS still has some remnants of the old master/slave or primary/secondary concept. The first server to mount each cluster file system becomes its primary; all other nodes in the cluster become secondaries. Applications access the user data in files directly from the server on which they are running. Each SFCFS node has its own intent

log. File system operations, such as allocating or deleting files, can originate from any node in the cluster.

Installing VxFS and enabling the cluster feature does not create a cluster file system configuration. File system clustering requires other Veritas products to enable communication services and provide storage resources. These products are packaged with VxFS in the Storage Foundation Cluster File System to provide a complete clustering environment.

See the *Veritas Storage Foundation Cluster File System Administrator's Guide*.

To be a cluster mount, a file system must be mounted using the `mount -o cluster` option. File systems mounted without the `-o cluster` option are termed local mounts.

SFCFS functionality is separately licensed.

Cross-platform data sharing

Cross-platform data sharing (CDS) allows data to be serially shared among heterogeneous systems where each system has direct access to the physical devices that hold the data. This feature can be used only in conjunction with Veritas Volume Manager (VxVM).

See the *Veritas Storage Foundation Cross-Platform Data Sharing Administrator's Guide*.

File Change Log

The VxFS File Change Log (FCL) tracks changes to files and directories in a file system. The File Change Log can be used by applications such as backup products, web crawlers, search and indexing engines, and replication software that typically scan an entire file system searching for modifications since a previous scan. FCL functionality is a separately licensed feature.

See [“About the File Change Log file”](#) on page 114.

Multi-volume support

The multi-volume support (MVS) feature allows several volumes to be represented by a single logical object. All I/O to and from an underlying logical volume is directed by way of volume sets. This feature can be used only in conjunction with VxVM. MVS functionality is a separately licensed feature.

See [“About multi-volume support”](#) on page 122.

Dynamic Storage Tiering

The Dynamic Storage Tiering (DST) option is built on multi-volume support technology. Using DST, you can map more than one volume to a single file system. You can then configure policies that automatically relocate files from one volume to another, or relocate files by running file relocation commands. Having multiple volumes lets you determine where files are located, which can improve performance for applications that access specific types of files. DST functionality is a separately licensed feature and is available with the `VRTSfppm` package.

See [“About Dynamic Storage Tiering”](#) on page 143.

Thin Reclamation of a file system

Storage is allocated from a Thin Storage LUN when files are created and written to a file system. This storage is not given back to the Thin Storage LUN when a file is deleted or the file size is shrunk. As such, the file system must perform the explicit task of releasing the free storage to the Thin Storage LUN. This is performed by the Storage Foundation Thin Reclamation feature. Thin Reclamation is only supported on VxFS file systems mounted on a VxVM volume.

Veritas File System performance enhancements

Traditional file systems employ block-based allocation schemes that provide adequate random access and latency for small files, but which limit throughput for larger files. As a result, they are less than optimal for commercial environments.

VxFS addresses this file system performance issue through an alternative allocation method and increased user control over allocation, I/O, and caching policies.

See [“Using Veritas File System”](#) on page 29.

VxFS provides the following performance enhancements:

- Data synchronous I/O
- Direct I/O and discovered direct I/O
- Support for files and file systems up to 256 terabytes
- Support for files up to 2 terabytes
- Enhanced I/O performance
- Caching advisories
- Enhanced directory features

- Explicit file alignment, extent size, and preallocation controls
- Tunable I/O parameters
- Tunable indirect data extent size
- Integration with VxVM™
- Support for large directories

Note: VxFS reduces the file lookup time in directories with an extremely large number of files.

About enhanced I/O performance

VxFS provides enhanced I/O performance by applying an aggressive I/O clustering policy, integrating with VxVM, and allowing application specific parameters to be set on a per-file system basis.

Enhanced I/O clustering

I/O clustering is a technique of grouping multiple I/O operations together for improved performance. VxFS I/O policies provide more aggressive clustering processes than other file systems and offer higher I/O throughput when using large files. The resulting performance is comparable to that provided by raw disk.

VxVM integration

VxFS interfaces with VxVM to determine the I/O characteristics of the underlying volume and perform I/O accordingly. VxFS also uses this information when using mkfs to perform proper allocation unit alignments for efficient I/O operations from the kernel. VxFS also uses this information when using mkfs to perform proper allocation unit alignments for efficient I/O operations from the kernel.

As part of VxFS/VxVM integration, VxVM exports a set of I/O parameters to achieve better I/O performance. This interface can enhance performance for different volume configurations such as RAID-5, striped, and mirrored volumes. Full stripe writes are important in a RAID-5 volume for strong I/O performance. VxFS uses these parameters to issue appropriate I/O requests to VxVM.

Application-specific parameters

You can also set application specific parameters on a per-file system basis to improve I/O performance.

- **Discovered Direct I/O**
All sizes above this value would be performed as direct I/O.
 - **Maximum Direct I/O Size**
This value defines the maximum size of a single direct I/O.
- See the `vxtunefs(1M)` and `tunefstab(4)` manual pages.

Using Veritas File System

There are three main methods to use, manage, modify, and tune VxFS:

- [Online system administration](#)
- [Application program interface](#)

Online system administration

VxFS provides command line interface (CLI) operations that are described throughout this guide and in manual pages.

VxFS allows you to run a number of administration tasks while the file system is online. Two of the more important tasks include:

- Defragmentation
- File system resizing

About defragmentation

Free resources are initially aligned and allocated to files in an order that provides optimal performance. On an active file system, the original order of free resources is lost over time as files are created, removed, and resized. The file system is spread farther along the disk, leaving unused gaps or fragments between areas that are in use. This process is known as fragmentation and leads to degraded performance because the file system has fewer options when assigning a free extent to a file (a group of contiguous data blocks).

VxFS provides the online administration utility `fsadm` to resolve the problem of fragmentation.

The `fsadm` utility defragments a mounted file system by performing the following actions:

- Removing unused space from directories
- Making all small files contiguous
- Consolidating free blocks for file system use

This utility can run on demand and should be scheduled regularly as a cron job.

About file system resizing

A file system is assigned a specific size as soon as it is created; the file system may become too small or too large as changes in file system usage take place over time.

VxFS is capable of increasing or decreasing the file system size while in use. Many competing file systems can not do this. The VxFS utility `fsadm` can expand or shrink a file system without unmounting the file system or interrupting user productivity. However, to expand a file system, the underlying device on which it is mounted must be expandable.

VxVM facilitates expansion using virtual disks that can be increased in size while in use. The VxFS and VxVM packages complement each other to provide online expansion capability. Use the `vxresize` command when resizing both the volume and the file system. The `vxresize` command guarantees that the file system shrinks or grows along with the volume. Do not use the `vxassist` and `fsadm_vxfs` commands for this purpose.

See the `vxresize(1M)` manual page.

See the *Veritas Volume Manager Administrator's Guide*.

Application program interface

Veritas File System Developer's Kit (SDK) provides developers with the information necessary to use the application programming interfaces (APIs) to modify and tune various features and components of File System.

See the *Veritas File System Programmer's Reference Guide*.

VxFS conforms to the System V Interface Definition (SVID) requirements and supports user access through the Network File System (NFS). Applications that require performance features not available with other file systems can take advantage of VxFS enhancements.

Expanded application facilities

VxFS provides API functions frequently associated with commercial applications that make it possible to perform the following actions:

- Preallocate space for a file
- Specify a fixed extent size for a file
- Bypass the system buffer cache for file I/O
- Specify the expected access pattern for a file

Because these functions are provided using VxFS-specific IOCTL system calls, most existing UNIX system applications do not use them. The `cp`, `cpio`, and `mv` utilities use the functions to preserve extent attributes and allocate space more efficiently. The current attributes of a file can be listed using the `getext` or `ls` command. The functions can also improve performance for custom applications. For portability reasons, these applications must check which file system type they are using before using these functions.

VxFS performance: creating, mounting, and tuning file systems

This chapter includes the following topics:

- [Creating a VxFS file system](#)
- [Mounting a VxFS file system](#)
- [Tuning the VxFS file system](#)
- [Monitoring free space](#)
- [Tuning I/O](#)

Creating a VxFS file system

When you create a file system with the `mkfs` command, you can select the following characteristics:

- [Block size](#)
- [Intent log size](#)

Block size

The unit of allocation in VxFS is a block. Unlike some other UNIX file systems, VxFS does not make use of block fragments for allocation because storage is allocated in extents that consist of one or more blocks.

You specify the block size when creating a file system by using the `mkfs -o bsize` option. The block size cannot be altered after the file system is created. The smallest available block size for VxFS is 1K. The default block size is 1024 bytes for file systems smaller than 1 TB, and 8192 bytes for file systems 1 TB or larger.

Choose a block size based on the type of application being run. For example, if there are many small files, a 1K block size may save space. For large file systems, with relatively few files, a larger block size is more appropriate. Larger block sizes use less disk space in file system overhead, but consume more space for files that are not a multiple of the block size. The easiest way to judge which block sizes provide the greatest system efficiency is to try representative system loads against various sizes and pick the fastest.

For 64-bit kernels, the block size and disk layout version determine the maximum size of the file system you can create.

See [“About disk layouts”](#) on page 255.

Intent log size

You specify the intent log size when creating a file system by using the `mkfs -o logsize` option. With the Version 6 or 7 disk layout, you can dynamically increase or decrease the intent log size using the `log` option of the `fsadm` command. The `mkfs` utility uses a default intent log size of 16 megabytes for disk layout Versions 4, 5, and 6. 64 megabytes is the default for Version 7. The default size is sufficient for most workloads. If the system is used as an NFS server or for intensive synchronous write workloads, performance may be improved using a larger log size. Version 7 is the default disk layout for VxFS 5.0.

With larger intent log sizes, recovery time is proportionately longer and the file system may consume more system resources (such as memory) during normal operation.

There are several system performance benchmark suites for which VxFS performs better with larger log sizes. As with block sizes, the best way to pick the log size is to try representative system loads against various sizes and pick the fastest.

Mounting a VxFS file system

In addition to the standard mount mode (`delaylog` mode), VxFS provides the following modes of operation:

- `log`
- `delaylog`

- `tmplog`
- `logsize`
- `nodatainlog`
- `blkclear`
- `mincache`
- `convosync`
- `ioerror`
- `largefiles|nolargefiles`
- `cio`
- `mntlock|mntunlock`
- `tranflush`

Caching behavior can be altered with the `mincache` option, and the behavior of `O_SYNC` and `D_SYNC` writes can be altered with the `convosync` option.

See the `fcntl(2)` manual page.

The `delaylog` and `tmplog` modes can significantly improve performance. The improvement over `log` mode is typically about 15 to 20 percent with `delaylog`; with `tmplog`, the improvement is even higher. Performance improvement varies, depending on the operations being performed and the workload. Read/write intensive loads should show less improvement, while file system structure intensive loads, such as `mkdir`, `create`, and `rename`, may show over 100 percent improvement. The best way to select a mode is to test representative system loads against the logging modes and compare the performance results.

Most of the modes can be used in combination. For example, a desktop machine might use both the `blkclear` and `mincache=closesync` modes.

See the `mount_vxfs(1M)` manual page.

The log mode

In `log` mode, all system calls other than `write(2)`, `writew(2)`, and `pwrite(2)` are guaranteed to be persistent after the system call returns to the application.

The `rename(2)` system call flushes the source file to disk to guarantee the persistence of the file data before renaming it. In both the `log` and `delaylog` modes, the `rename` is also guaranteed to be persistent when the system call returns. This benefits shell scripts and programs that try to update a file atomically by

writing the new file contents to a temporary file and then renaming it on top of the target file.

The delaylog mode

The default logging mode is `delaylog`. In `delaylog` mode, the effects of most system calls other than `write(2)`, `writew(2)`, and `pwrite(2)` are guaranteed to be persistent approximately 15 to 20 seconds after the system call returns to the application. Contrast this with the behavior of most other file systems in which most system calls are not persistent until approximately 30 seconds or more after the call has returned. Fast file system recovery works with this mode.

The `rename(2)` system call flushes the source file to disk to guarantee the persistence of the file data before renaming it. In the `log` and `delaylog` modes, the `rename` is also guaranteed to be persistent when the system call returns. This benefits shell scripts and programs that try to update a file atomically by writing the new file contents to a temporary file and then renaming it on top of the target file.

The tmplog mode

In `tmplog` mode, the effects of system calls have persistence guarantees that are similar to those in `delaylog` mode. In addition, enhanced flushing of delayed extending writes is disabled, which results in better performance but increases the chances of data being lost or uninitialized data appearing in a file that was being actively written at the time of a system failure. This mode is only recommended for temporary file systems. Fast file system recovery works with this mode.

Note: The term "effects of system calls" refers to changes to file system data and metadata caused by the system call, excluding changes to `st_atime`.

See the `stat(2)` manual page.

Persistence guarantees

In all logging modes, VxFS is fully POSIX compliant. The effects of the `fsync(2)` and `fdatasync(2)` system calls are guaranteed to be persistent after the calls return. The persistence guarantees for data or metadata modified by `write(2)`, `writew(2)`, or `pwrite(2)` are not affected by the logging mount options. The effects of these system calls are guaranteed to be persistent only if the `O_SYNC`, `O_DSYNC`, `VX_DSYNC`, or `VX_DIRECT` flag, as modified by the `convosync=` mount option, has been specified for the file descriptor.

The behavior of NFS servers on a VxFS file system is unaffected by the `log` and `tmplog` mount options, but not `delaylog`. In all cases except for `tmplog`, VxFS complies with the persistency requirements of the NFS v2 and NFS v3 standard. Unless a UNIX application has been developed specifically for the VxFS file system in `log` mode, it expects the persistence guarantees offered by most other file systems and experiences improved robustness when used with a VxFS file system mounted in `delaylog` mode. Applications that expect better persistence guarantees than that offered by most other file systems can benefit from the `log`, `mincache=`, and `closesync` mount options. However, most commercially available applications work well with the default VxFS mount options, including the `delaylog` mode.

The logiosize mode

The `logiosize=size` option enhances the performance of storage devices that employ a read-modify-write feature. If you specify `logiosize` when you mount a file system, VxFS writes the intent log in the least *size* bytes or a multiple of *size* bytes to obtain the maximum performance from such devices.

See the `mount_vxfs(1M)` manual page.

The values for *size* can be 1024, 2048, or 4096.

The nodatainlog mode

Use the `nodatainlog` mode on systems with disks that do not support bad block revectoring. Usually, a VxFS file system uses the intent log for synchronous writes. The inode update and the data are both logged in the transaction, so a synchronous write only requires one disk write instead of two. When the synchronous write returns to the application, the file system has told the application that the data is already written. If a disk error causes the metadata update to fail, then the file must be marked bad and the entire file is lost.

If a disk supports bad block revectoring, then a failure on the data update is unlikely, so logging synchronous writes should be allowed. If the disk does not support bad block revectoring, then a failure is more likely, so the `nodatainlog` mode should be used.

A `nodatainlog` mode file system is approximately 50 percent slower than a standard mode VxFS file system for synchronous writes. Other operations are not affected.

The blkclear mode

The `blkclear` mode is used in increased data security environments. The `blkclear` mode guarantees that uninitialized storage never appears in files. The increased integrity is provided by clearing extents on disk when they are allocated within a file. This mode does not affect extending writes. A `blkclear` mode file system is approximately 10 percent slower than a standard mode VxFS file system, depending on the workload.

The mincache mode

The mincache mode has the following suboptions:

- `mincache=closesync`
- `mincache=direct`
- `mincache=dsync`
- `mincache=unbuffered`
- `mincache=tmpcache`

The `mincache=closesync` mode is useful in desktop environments where users are likely to shut off the power on the machine without halting it first. In this mode, any changes to the file are flushed to disk when the file is closed.

To improve performance, most file systems do not synchronously update data and inode changes to disk. If the system crashes, files that have been updated within the past minute are in danger of losing data. With the `mincache=closesync` mode, if the system crashes or is switched off, only open files can lose data. A `mincache=closesync` mode file system could be approximately 15 percent slower than a standard mode VxFS file system, depending on the workload.

The following describes where to use the mincache modes:

- The `mincache=direct`, `mincache=unbuffered`, and `mincache=dsync` modes are used in environments where applications have reliability problems caused by the kernel buffering of I/O and delayed flushing of non-synchronous I/O.
- The `mincache=direct` and `mincache=unbuffered` modes guarantee that all non-synchronous I/O requests to files are handled as if the `VX_DIRECT` or `VX_UNBUFFERED` caching advisories had been specified.
- The `mincache=dsync` mode guarantees that all non-synchronous I/O requests to files are handled as if the `VX_DSYNC` caching advisory had been specified. Refer to the `vxfsio(7)` manual page for explanations of `VX_DIRECT`, `VX_UNBUFFERED`, and `VX_DSYNC`, as well as for the requirements for direct I/O.

- The `mincache=direct`, `mincache=unbuffered`, and `mincache=dsync` modes also flush file data on close as `mincache=closesync` does.

Because the `mincache=direct`, `mincache=unbuffered`, and `mincache=dsync` modes change non-synchronous I/O to synchronous I/O, throughput can substantially degrade for small to medium size files with most applications. Since the `VX_DIRECT` and `VX_UNBUFFERED` advisories do not allow any caching of data, applications that normally benefit from caching for reads usually experience less degradation with the `mincache=dsync` mode. `mincache=direct` and `mincache=unbuffered` require significantly less CPU time than buffered I/O.

If performance is more important than data integrity, you can use the `mincache=tmpcache` mode. The `mincache=tmpcache` mode disables special delayed extending write handling, trading off less integrity for better performance. Unlike the other `mincache` modes, `tmpcache` does not flush the file to disk the file is closed. When the `mincache=tmpcache` option is used, bad data can appear in a file that was being extended when a crash occurred.

The convosync mode

The `convosync` (convert `osync`) mode has the following suboptions:

- `convosync=closesync`

Note: The `convosync=closesync` mode converts synchronous and data synchronous writes to non-synchronous writes and flushes the changes to the file to disk when the file is closed.

- `convosync=delay`
- `convosync=direct`
- `convosync=dsync`

Note: The `convosync=dsync` option violates POSIX guarantees for synchronous I/O.

- `convosync=unbuffered`

The `convosync=delay` mode causes synchronous and data synchronous writes to be delayed rather than to take effect immediately. No special action is performed when closing a file. This option effectively cancels any data integrity guarantees normally provided by opening a file with `O_SYNC`.

See the `open(2)`, `fcntl(2)`, and `vxfsio(7)` manual pages.

Warning: Be very careful when using the `convosync=closesync` or `convosync=delay` mode because they actually change synchronous I/O into non-synchronous I/O. Applications that use synchronous I/O for data reliability may fail if the system crashes and synchronously written data is lost.

The `convosync=dsync` mode converts synchronous writes to data synchronous writes.

As with `closesync`, the `direct`, `unbuffered`, and `dsync` modes flush changes to the file to disk when it is closed. These modes can be used to speed up applications that use synchronous I/O. Many applications that are concerned with data integrity specify the `O_SYNC` `fcntl` in order to write the file data synchronously. However, this has the undesirable side effect of updating inode times and therefore slowing down performance. The `convosync=dsync`, `convosync=unbuffered`, and `convosync=direct` modes alleviate this problem by allowing applications to take advantage of synchronous writes without modifying inode times as well.

Before using `convosync=dsync`, `convosync=unbuffered`, or `convosync=direct`, make sure that all applications that use the file system do not require synchronous inode time updates for `O_SYNC` writes.

The `ioerror` mode

This mode sets the policy for handling I/O errors on a mounted file system. I/O errors can occur while reading or writing file data or metadata. The file system can respond to these I/O errors either by halting or by gradually degrading. The `ioerror` option provides five policies that determine how the file system responds to the various errors. All policies limit data corruption, either by stopping the file system or by marking a corrupted inode as bad.

The policies are the following:

- `disable`
- `nodisable`
- `wdisable`
- `mwdisable`
- `mdisable`

The disable policy

If `disable` is selected, VxFS disables the file system after detecting any I/O error. You must then unmount the file system and correct the condition causing the I/O error. After the problem is repaired, run `fsck` and mount the file system again. In most cases, replay `fsck` to repair the file system. A full `fsck` is required only in cases of structural damage to the file system's metadata. Select `disable` in environments where the underlying storage is redundant, such as RAID-5 or mirrored disks.

The nodisable policy

If `nodisable` is selected, when VxFS detects an I/O error, it sets the appropriate error flags to contain the error, but continues running. Note that the degraded condition indicates possible data or metadata corruption, not the overall performance of the file system.

For file data read and write errors, VxFS sets the `VX_DATAIOERR` flag in the super-block. For metadata read errors, VxFS sets the `VX_FULLFSCK` flag in the super-block. For metadata write errors, VxFS sets the `VX_FULLFSCK` and `VX_METAIOERR` flags in the super-block and may mark associated metadata as bad on disk. VxFS then prints the appropriate error messages to the console.

See [“File system response to problems”](#) on page 209.

You should stop the file system as soon as possible and repair the condition causing the I/O error. After the problem is repaired, run `fsck` and mount the file system again. Select `nodisable` if you want to implement the policy that most closely resembles the error handling policy of the previous VxFS release.

The wdisable and mwdisable policies

If `wdisable` (write disable) or `mwdisable` (metadata-write disable) is selected, the file system is disabled or degraded, depending on the type of error encountered. Select `wdisable` or `mwdisable` for environments where read errors are more likely to persist than write errors, such as when using non-redundant storage. `mwdisable` is the default `ioerror` mount option for local mounts.

See the `mount_vxfs(1M)` manual page.

The mdisable policy

If `mdisable` (metadata disable) is selected, the file system is disabled if a metadata read or write fails. However, the file system continues to operate if the failure is

confined to data extents. `mdisable` is the default `ioerror` mount option for cluster mounts.

The `largefiles|nolargefiles` option

The section includes the following topics :

- [Creating a file system with large files](#)
- [Mounting a file system with large files](#)
- [Managing a file system with large files](#)

VxFS supports files larger than 2 gigabytes. The maximum file size that can be created is 2 terabytes.

Note: Applications and utilities such as backup may experience problems if they are not aware of large files. In such a case, create your file system without large file capability.

Creating a file system with large files

To create a file system with a file capability:

```
# mkfs -F vxfs -o largefiles special_device size
```

Specifying `largefiles` sets the `largefiles` flag. This lets the file system to hold files that are two gigabytes or larger. This is the default option.

To clear the flag and prevent large files from being created:

```
# mkfs -F vxfs -o nolargefiles special_device size
```

The `largefiles` flag is persistent and stored on disk.

Mounting a file system with large files

If a mount succeeds and `nolargefiles` is specified, the file system cannot contain or create any large files. If a mount succeeds and `largefiles` is specified, the file system may contain and create large files.

The `mount` command fails if the specified `largefiles|nolargefiles` option does not match the on-disk flag.

Because the `mount` command defaults to match the current setting of the on-disk flag if specified without the `largefiles` or `nolargefiles` option, the best practice

is not to specify either option. After a file system is mounted, you can use the `fsadm` utility to change the large files option.

Managing a file system with large files

Managing a file system with large files includes the following tasks:

- Determining the current status of the large files flag
- Switching capabilities on a mounted file system
- Switching capabilities on an unmounted file system

To determine the current status of the `largefiles` flag, type either of the following commands:

```
# mkfs -F vxfs -m special_device
# fsadm -F vxfs mount_point | special_device
```

To switch capabilities on a mounted file system:

```
# fsadm -F vxfs -o [no]largefiles mount_point
```

To switch capabilities on an unmounted file system:

```
# fsadm -F vxfs -o [no]largefiles special_device
```

You cannot change a file system to `nolargefiles` if it contains large files.

See the `mount_vxfs(1M)`, `fsadm_vxfs(1M)`, and `mkfs_vxfs(1M)` manual pages.

The cio option

The `cio` (Concurrent I/O) option specifies the file system to be mounted for concurrent readers and writers. Concurrent I/O is a licensed feature of VxFS. If `cio` is specified, but the feature is not licensed, the `mount` command prints an error message and terminates the operation without mounting the file system. The `cio` option cannot be disabled through a remount. To disable the `cio` option, the file system must be unmounted and mounted again without the `cio` option.

The mntlock|mntunlock option

The `mntlock` option prevents a file system from being unmounted by an application. This option is useful for applications that do not want the file systems that the applications are monitoring to be improperly unmounted by other applications or administrators.

The `mntunlock` option of the `vxumount` command reverses the `mntlock` option if you previously locked the file system.

Combining mount command options

Although mount options can be combined arbitrarily, some combinations do not make sense. The following examples provide some common and reasonable mount option combinations.

To mount a desktop file system using options:

```
# mount -F vxfs -o log,mincache=closesync /dev/dsk/c1t3d0 /mnt
```

This guarantees that when a file is closed, its data is synchronized to disk and cannot be lost. Thus, after an application has exited and its files are closed, no data is lost even if the system is immediately turned off.

To mount a temporary file system or to restore from backup:

```
# mount -F vxfs -o tmplog,convosync=delay,mincache=tmpcache \  
/dev/dsk/c1t3d0 /mnt
```

This combination might be used for a temporary file system where performance is more important than absolute data integrity. Any `O_SYNC` writes are performed as delayed writes and delayed extending writes are not handled. This could result in a file that contains corrupted data if the system crashes. Any file written 30 seconds or so before a crash may contain corrupted data or be missing if this mount combination is in effect. However, such a file system does significantly less disk writes than a log file system, and should have significantly better performance, depending on the application.

To mount a file system for synchronous writes:

```
# mount -F vxfs -o log,convosync=dsync /dev/dsk/c1t3d0 /mnt
```

This combination can be used to improve the performance of applications that perform `O_SYNC` writes, but only require data synchronous write semantics. Performance can be significantly improved if the file system is mounted using `convosync=dsync` without any loss of data integrity.

Tuning the VxFS file system

This section describes the following kernel tunable parameters in VxFS:

- [Tuning inode table size](#)
- [\[Unresolved xref\]](#)

- [VxFS buffer cache high water mark](#)
- [Number of links to a file](#)
- [VxFS inode free time lag](#)

Tuning inode table size

VxFS caches inodes in an inode table. There is a dynamic tunable in VxFS called `vx_ninode` that determines the number of entries in the inode table. You can dynamically change the value of `vx_ninode` by using the `sam` or `kctune` commands.

See the `sam(1M)` and `kctune(1M)` manual pages.

This value is used to determine the number of entries in the VxFS inode table. By default, `vx_ninode` initializes at zero; the file system then computes a value based on the system memory size.

A VxFS file system can also obtain the value of `vx_ninode` from the system configuration file used for making the HP-UX kernel (`/stand/system` for example). To change the computed value of `vx_ninode`, you can add an entry to the `system` configuration file. For example:

```
tunable vx_ninode 1000000
```

This sets the inode table size to 1,000,000 inodes after making a new HP-UX kernel using `mk_kernel`.

Increasing the value of `vx_ninode` increases the inode table size immediately, allowing a higher number of inodes to be cached. Decreasing the value of `vx_ninode` decreases the inode table size to the specified value. After the tunable is decreased, VxFS attempts to free excess cached objects so that the resulting number of inodes in the table is less than or equal to the specified value of `vx_ninode`. If this attempt fails, the value of the `vx_ninode` tunable is not changed. In such a case, the `kctune` command can be specified with the `-h` option so that the new value of `vx_ninode` takes effect after a system reboot.

Be careful when changing the value of `vx_ninode`, as the value can affect file system performance. Typically, the default value determined by VxFS based on the amount of system memory ensures good system performance across a wide range of applications. However, if it is determined that the default value is not suitable, `vx_ninode` can be set to an appropriate value based on the expected file system usage. The `vxfsstat` command can be used to monitor inode cache usage and statistics to determine the optimum value of `vx_ninode` for the system.

Changing the value of a tunable does not resize the internal hash tables and structures of the caches. These sizes are determined at system boot up based on

either the system memory size, which is the default, or the value of the tunable if explicitly set, whichever is larger. Thus, dynamically increasing the tunable to a value that is more than two times either the default value or the user-defined value, if larger, may cause performance degradation unless the system is rebooted.

Examples of changing the vx_inode tunable value

The following are examples of changing the `vx_inode` tunable value.

Reporting the current value of vx_inode

```
# kctune vx_inode
```

This command displays the current value of `vx_inode`.

Setting vx_inode

```
# kctune -s vx_inode=10000
```

This command sets `vx_inode` to 10000, the specified value.

Restoring vx_inode to its default value

```
# kctune -s vx_inode=
```

This command restores `vx_inode` to its default value by clearing the user-specified value. The default value is the value determined by VxFS to be optimal based on the amount of system memory, which is used if `vx_inode` is not explicitly set.

Delaying a change to vx_inode until after a reboot

```
# kctune -h -s vx_inode=10000
```

If the `-h` option is specified, the specified value for `vx_inode` does not take effect until after a system reboot.

VxFS buffer cache high water mark

VxFS maintains its own buffer cache in the kernel for frequently accessed file system metadata. This cache is different from the HP-UX kernel buffer cache that caches file data. The `vx_bc_bufhwm` dynamic, global, tunable parameter lets you change the VxFS buffer cache high water mark, that is, the maximum amount of memory that can be used to cache VxFS metadata.

The initial value of `vx_bc_bufhwm` is zero. When the operating system reboots, VxFS sets the value of `vx_bc_bufhwm` based on the amount of system memory. You can explicitly reset the value of `vx_bc_bufhwm` by changing the value of `vxfs_bc_bufhwm` using the `sam` or `kctune` commands.

See the `sam(1M)` and `kctune(1M)` manual pages.

You can also set the value by adding an entry to the system configuration file. For example, the following entry:

```
vxfs_bc_bufhwm vx_bc_bufhwm 300000
```

sets the high water mark to 300 megabytes. The change takes effect after you rebuild the HP-UX kernel using the `mk_kernel` command. You specify the `vx_bc_bufhwm` tunable in units of kilobytes. The minimum value is 6144.

Increasing the value of `vx_bc_bufhwm` increases the VxFS buffer cache immediately, allowing a greater amount of memory to be used to cache VxFS metadata. Decreasing the value of `vx_bc_bufhwm` decreases the VxFS buffer cache to the specified value. This frees memory such that the amount of memory used for buffer cache is lower than the specified value of `vx_bc_bufhwm`.

Typically, the default value computed by VxFS based on the amount of system memory ensures good system performance across a wide range of applications. For application loads that cause frequent file system metadata changes on the system (for example, a high rate of file creation or deletion, or accessing large directories), changing the value of `vx_bc_bufhwm` may improve performance.

You can use the `vxfsstat` command to monitor buffer cache statistics and inode cache usage.

See the `vxfsstat(1M)` manual page.

Number of links to a file

In VxFS, the number of possible links to a file is determined by the `vx_maxlink` global tunable. The default value of `vx_maxlink` is 32767, the maximum value is 65535. This is a static tunable.

You can set the value of `vx_maxlink` using the `sam` or `kctune` commands.

See the `sam(1M)` and `kctune(1M)` manual pages.

You can also add an entry to the system configuration file as shown in the following example:

```
vxfs_maxlink vx_maxlink 40000
```

This sets the value of `vx_maxlink` to 40,000 links.

VxFS inode free time lag

In VxFS, an inode is put on a freelist if it is not being used. The memory space for this unused inode can be freed if it stays on the freelist for a specified amount of time. The `vx_ifree_timelag` tunable specifies the minimum amount of time an unused inode spends on a freelist before its memory space is freed.

The `vx_ifree_timelag` tunable is dynamic. Any changes to `vx_ifree_timelag` take effect immediately.

The default value of `vx_ifree_timelag` is 0. By setting `vx_ifree_timelag` to 0, the inode free time lag is autotuned to 1800 seconds. Specifying negative one (-1) stops the freeing of inode space; no further inode allocations are freed until the value is changed back to a value other than negative one.

You can change the value of `vx_ifree_timelag` using the `sam` or `kctune` commands.

See the `sam(1M)` and `kctune(1M)` manual pages.

You can also add an entry to the system configuration file. The following example changes the value of `vx_ifree_timelag` to 2400 seconds:

```
# kctune -s vxfs_ifree_timelag=2400
```

Note: The default value of `vx_ifree_timelag` typically provides optimal VxFS performance. Be careful when adjusting the tunable because incorrect tuning can adversely affect system performance.

Monitoring free space

In general, VxFS works best if the percentage of free space in the file system does not get below 10 percent. This is because file systems with 10 percent or more free space have less fragmentation and better extent allocation. Regular use of the `df` command to monitor free space is desirable.

See the `df_vxfs(1M)` manual page.

Full file systems may have an adverse effect on file system performance. Full file systems should therefore have some files removed, or should be expanded.

See the `fsadm_vxfs(1M)` manual page.

Monitoring fragmentation

Fragmentation reduces performance and availability. Regular use of `fsadm`'s fragmentation reporting and reorganization facilities is therefore advisable.

The easiest way to ensure that fragmentation does not become a problem is to schedule regular defragmentation runs using the `cron` command.

Defragmentation scheduling should range from weekly (for frequently used file systems) to monthly (for infrequently used file systems). Extent fragmentation should be monitored with `fsadm` command.

To determine the degree of fragmentation, use the following factors:

- Percentage of free space in extents of less than 8 blocks in length
- Percentage of free space in extents of less than 64 blocks in length
- Percentage of free space in extents of length 64 blocks or greater

An unfragmented file system has the following characteristics:

- Less than 1 percent of free space in extents of less than 8 blocks in length
- Less than 5 percent of free space in extents of less than 64 blocks in length
- More than 5 percent of the total file system size available as free extents in lengths of 64 or more blocks

A badly fragmented file system has one or more of the following characteristics:

- Greater than 5 percent of free space in extents of less than 8 blocks in length
- More than 50 percent of free space in extents of less than 64 blocks in length
- Less than 5 percent of the total file system size available as free extents in lengths of 64 or more blocks

The optimal period for scheduling of extent reorganization runs can be determined by choosing a reasonable interval, scheduling `fsadm` runs at the initial interval, and running the extent fragmentation report feature of `fsadm` before and after the reorganization.

The “before” result is the degree of fragmentation prior to the reorganization. If the degree of fragmentation is approaching the figures for bad fragmentation, reduce the interval between `fsadm` runs. If the degree of fragmentation is low, increase the interval between `fsadm` runs.

The “after” result is an indication of how well the reorganizer has performed. The degree of fragmentation should be close to the characteristics of an unfragmented file system. If not, it may be a good idea to resize the file system; full file systems tend to fragment and are difficult to defragment. It is also possible that the

reorganization is not being performed at a time during which the file system in question is relatively idle.

Directory reorganization is not nearly as critical as extent reorganization, but regular directory reorganization improves performance. It is advisable to schedule directory reorganization for file systems when the extent reorganization is scheduled. The following is a sample script that is run periodically at 3:00 A.M. from `cron` for a number of file systems:

```
outfile=/usr/spool/fsadm/out.`/bin/date +%m%d`
for i in /home /home2 /project /db
do
    /bin/echo "Reorganizing $i"
    /bin/timex fsadm -F vxfs -e -E -s $i
    /bin/timex fsadm -F vxfs -s -d -D $i
done > $outfile 2>&1
```

Thin Reclamation

Veritas File System (VxFS) supports reclamation of free storage on a Thin Storage LUN. Free storage is reclaimed using the `fsadm` command or the `vxfs_ts_reclaim` API. You can perform the default reclamation or aggressive reclamation. If you used a file system for a long time and must perform reclamation on the file system, Symantec recommends that you run aggressive reclamation. Aggressive reclamation compacts the allocated blocks, which creates larger free blocks that can potentially be reclaimed.

See the `fsadm_vxfs(1M)` and `vxfs_ts_reclaim(3)` manual pages.

Thin Reclamation is only supported on file systems mounted on a VxVM volume.

The following example performs reclamation of free storage to the Thin Storage LUN on a VxFS file system mounted at `/mnt1`:

```
# fsadm -R /mnt1
```

Veritas File System also supports reclamation of a portion of the file system using the `vxfs_ts_reclaim()` API.

See the *Veritas File System Programmer's Reference Guide*.

Note: Thin Reclamation is a slow process and may take several hours to complete, depending on the file system size. Thin Reclamation is not guaranteed to reclaim 100% of the free space.

You can track the progress of the Thin Reclamation process by using the `vxtask list` command when using the Veritas Volume Manager (VxVM) command `vxdisk reclaim`.

See the `vxtask(1M)` and `vxdisk(1M)` manual pages.

You can administer Thin Reclamation using VxVM commands.

See the *Veritas Volume Manager Administrator's Guide*.

Tuning I/O

The performance of a file system can be enhanced by a suitable choice of I/O sizes and proper alignment of the I/O requests based on the requirements of the underlying special device. VxFS provides tools to tune the file systems.

Note: The following tunables and the techniques work on a per file system basis. Use them judiciously based on the underlying device properties and characteristics of the applications that use the file system.

Tuning VxFS I/O parameters

VxFS provides a set of tunable I/O parameters that control some of its behavior. These I/O parameters are useful to help the file system adjust to striped or RAID-5 volumes that could yield performance superior to a single disk. Typically, data streaming applications that access large files see the largest benefit from tuning the file system.

VxVM queries

VxVM receives the following queries during configuration:

- The file system queries VxVM to determine the geometry of the underlying volume and automatically sets the I/O parameters.

Note: When using file systems in multiple volume sets, VxFS sets the VxFS tunables based on the geometry of the first component volume (volume 0) in the volume set.

- The `mount` command queries VxVM when the file system is mounted and downloads the I/O parameters.

If the default parameters are not acceptable or the file system is being used without VxVM, then the `/etc/vx/tunefstab` file can be used to set values for I/O parameters. The `mount` command reads the `/etc/vx/tunefstab` file and downloads any parameters specified for a file system. The `tunefstab` file overrides any values obtained from VxVM. While the file system is mounted, any I/O parameters can be changed using the `vxtunefs` command which can have tunables specified on the command line or can read them from the `/etc/vx/tunefstab` file.

See the `vxtunefs(1M)` and `tunefstab(4)` manual pages.

The `vxtunefs` command can be used to print the current values of the I/O parameters.

To print the values, type the following command:

```
# vxtunefs -p mount_point
```

The following is an example `tunefstab` file:

```
/dev/vx/dsk/userdg/netbackup
read_pref_io=128k,write_pref_io=128k,read_nstream=4,write_nstream=4
/dev/vx/dsk/userdg/metasave
read_pref_io=128k,write_pref_io=128k,read_nstream=4,write_nstream=4
/dev/vx/dsk/userdg/solbuild
read_pref_io=64k,write_pref_io=64k,read_nstream=4,write_nstream=4
/dev/vx/dsk/userdg/solrelease
read_pref_io=64k,write_pref_io=64k,read_nstream=4,write_nstream=4
/dev/vx/dsk/userdg/solpatch
read_pref_io=128k,write_pref_io=128k,read_nstream=4,write_nstream=4
```

Tunable I/O parameters

Table 2-1 provides a list and description of these parameters.

Table 2-1 Tunable VxFS I/O parameters

Parameter	Description
<code>read_pref_io</code>	The preferred read request size. The file system uses this in conjunction with the <code>read_nstream</code> value to determine how much data to read ahead. The default value is 64K.

Table 2-1 Tunable VxFS I/O parameters (*continued*)

Parameter	Description
<code>write_pref_io</code>	The preferred write request size. The file system uses this in conjunction with the <code>write_nstream</code> value to determine how to do flush behind on writes. The default value is 64K.
<code>read_nstream</code>	The number of parallel read requests of size <code>read_pref_io</code> to have outstanding at one time. The file system uses the product of <code>read_nstream</code> multiplied by <code>read_pref_io</code> to determine its read ahead size. The default value for <code>read_nstream</code> is 1.
<code>write_nstream</code>	The number of parallel write requests of size <code>write_pref_io</code> to have outstanding at one time. The file system uses the product of <code>write_nstream</code> multiplied by <code>write_pref_io</code> to determine when to do flush behind on writes. The default value for <code>write_nstream</code> is 1.
<code>discovered_direct_iosz</code>	Any file I/O requests larger than <code>discovered_direct_iosz</code> are handled as discovered direct I/O. A discovered direct I/O is unbuffered similar to direct I/O, but it does not require a synchronous commit of the inode when the file is extended or blocks are allocated. For larger I/O requests, the CPU time for copying the data into the page cache and the cost of using memory to buffer the I/O data becomes more expensive than the cost of doing the disk I/O. For these I/O requests, using discovered direct I/O is more efficient than regular I/O. The default value of this parameter is 256K.

Table 2-1 Tunable VxFS I/O parameters (continued)

Parameter	Description
<code>fcl_keeptime</code>	<p>Specifies the minimum amount of time, in seconds, that the VxFS File Change Log (FCL) keeps records in the log. When the oldest 8K block of FCL records have been kept longer than the value of <code>fcl_keeptime</code>, they are purged from the FCL and the extents nearest to the beginning of the FCL file are freed. This process is referred to as "punching a hole." Holes are punched in the FCL file in 8K chunks.</p> <p>If the <code>fcl_maxalloc</code> parameter is set, records are purged from the FCL if the amount of space allocated to the FCL exceeds <code>fcl_maxalloc</code>, even if the elapsed time the records have been in the log is less than the value of <code>fcl_keeptime</code>. If the file system runs out of space before <code>fcl_keeptime</code> is reached, the FCL is deactivated.</p> <p>Either or both of the <code>fcl_keeptime</code> or <code>fcl_maxalloc</code> parameters must be set before the File Change Log can be activated. <code>fcl_keeptime</code> does not apply to disk layout Versions 1 through 5.</p>
<code>fcl_maxalloc</code>	<p>Specifies the maximum amount of space that can be allocated to the VxFS File Change Log (FCL). The FCL file is a sparse file that grows as changes occur in the file system. When the space allocated to the FCL file reaches the <code>fcl_maxalloc</code> value, the oldest FCL records are purged from the FCL and the extents nearest to the beginning of the FCL file are freed. This process is referred to as "punching a hole." Holes are punched in the FCL file in 8K chunks. If the file system runs out of space before <code>fcl_maxalloc</code> is reached, the FCL is deactivated.</p> <p>The minimum value of <code>fcl_maxalloc</code> is 4 MB. The default value is <code>fs_size/33</code>.</p> <p>Either or both of the <code>fcl_maxalloc</code> or <code>fcl_keeptime</code> parameters must be set before the File Change Log can be activated. <code>fcl_maxalloc</code> does not apply to disk layout Versions 1 through 5.</p>

Table 2-1 Tunable VxFS I/O parameters (*continued*)

Parameter	Description
<code>fcl_winterval</code>	<p>Specifies the time, in seconds, that must elapse before the VxFS File Change Log (FCL) records a data overwrite, data extending write, or data truncate for a file. The ability to limit the number of repetitive FCL records for continuous writes to the same file is important for file system performance and for applications processing the FCL. <code>fcl_winterval</code> is best set to an interval less than the shortest interval between reads of the FCL by any application. This way all applications using the FCL can be assured of finding at least one FCL record for any file experiencing continuous data changes.</p> <p><code>fcl_winterval</code> is enforced for all files in the file system. Each file maintains its own time stamps, and the elapsed time between FCL records is per file. This elapsed time can be overridden using the VxFS FCL sync public API.</p> <p>See the <code>vxfs_fcl_sync(3)</code> manual page.</p> <p><code>fcl_winterval</code> does not apply to disk layout Versions 1 through 5.</p>
<code>hsm_write_prealloc</code>	<p>For a file managed by a hierarchical storage management (HSM) application, <code>hsm_write_prealloc</code> preallocates disk blocks before data is migrated back into the file system. An HSM application usually migrates the data back through a series of writes to the file, each of which allocates a few blocks. By setting <code>hsm_write_prealloc</code> (<code>hsm_write_prealloc=1</code>), a sufficient number of disk blocks are allocated on the first write to the empty file so that no disk block allocation is required for subsequent writes. This improves the write performance during migration.</p> <p>The <code>hsm_write_prealloc</code> parameter is implemented outside of the DMAPI specification, and its usage has limitations depending on how the space within an HSM-controlled file is managed. It is advisable to use <code>hsm_write_prealloc</code> only when recommended by the HSM application controlling the file system.</p>

Table 2-1 Tunable VxFS I/O parameters (continued)

Parameter	Description
<code>initial_extent_size</code>	<p>Changes the default initial extent size. VxFS determines, based on the first write to a new file, the size of the first extent to be allocated to the file. Normally the first extent is the smallest power of 2 that is larger than the size of the first write. If that power of 2 is less than 8K, the first extent allocated is 8K. After the initial extent, the file system increases the size of subsequent extents with each allocation.</p> <p>See <code>max_seqio_extent_size</code>.</p> <p>Since most applications write to files using a buffer size of 8K or less, the increasing extents start doubling from a small initial extent. <code>initial_extent_size</code> can change the default initial extent size to be larger, so the doubling policy starts from a much larger initial size and the file system does not allocate a set of small extents at the start of file. Use this parameter only on file systems that have a very large average file size. On these file systems it results in fewer extents per file and less fragmentation. <code>initial_extent_size</code> is measured in file system blocks.</p>
<code>inode_aging_count</code>	<p>Specifies the maximum number of inodes to place on an inode aging list. Inode aging is used in conjunction with file system Storage Checkpoints to allow quick restoration of large, recently deleted files. The aging list is maintained in first-in-first-out (fifo) order up to maximum number of inodes specified by <code>inode_aging_count</code>. As newer inodes are placed on the list, older inodes are removed to complete their aging process. For best performance, it is advisable to age only a limited number of larger files before completion of the removal process. The default maximum number of inodes to age is 2048.</p>
<code>max_buf_data_size</code>	<p>The maximum buffer size allocated for file data; either 8K bytes or 64K bytes. Use the larger value for workloads where large reads/writes are performed sequentially. Use the smaller value on workloads where the I/O is random or is done in small chunks. 8K bytes is the default value.</p>

Table 2-1 Tunable VxFS I/O parameters (*continued*)

Parameter	Description
<code>inode_aging_size</code>	Specifies the minimum size to qualify a deleted inode for inode aging. Inode aging is used in conjunction with file system Storage Checkpoints to allow quick restoration of large, recently deleted files. For best performance, it is advisable to age only a limited number of larger files before completion of the removal process. Setting the size too low can push larger file inodes out of the aging queue to make room for newly removed smaller file inodes.
<code>max_direct_iosz</code>	The maximum size of a direct I/O request that are issued by the file system. If a larger I/O request comes in, then it is broken up into <code>max_direct_iosz</code> chunks. This parameter defines how much memory an I/O request can lock at once, so it should not be set to more than 20 percent of memory.
<code>max_diskq</code>	Limits the maximum disk queue generated by a single file. When the file system is flushing data for a file and the number of buffers being flushed exceeds <code>max_diskq</code> , processes are blocked until the amount of data being flushed decreases. Although this does not limit the actual disk queue, it prevents flushing processes from making the system unresponsive. The default value is 1 MB.
<code>max_seqio_extent_size</code>	Increases or decreases the maximum size of an extent. When the file system is following its default allocation policy for sequential writes to a file, it allocates an initial extent which is large enough for the first write to the file. When additional extents are allocated, they are progressively larger because the algorithm tries to double the size of the file with each new extent. As such, each extent can hold several writes worth of data. This is done to reduce the total number of extents in anticipation of continued sequential writes. When the file stops being written, any unused space is freed for other files to use. Normally, this allocation stops increasing the size of extents at 262144 blocks, which prevents one file from holding too much unused space. <code>max_seqio_extent_size</code> is measured in file system blocks. The default and minimum value of is 2048 blocks.

Table 2-1 Tunable VxFS I/O parameters (continued)

Parameter	Description
default_indir_size	<p>On VxFS, files can have up to ten direct extents of variable size stored in the inode. After these extents are used up, the file must use indirect extents which are a fixed size that is set when the file first uses indirect extents. These indirect extents are 8K by default. The file system does not use larger indirect extents because it must fail a write and return <code>ENOSPC</code> if there are no extents available that are the indirect extent size. For file systems with many large files, the 8K indirect extent size is too small. The files that get into indirect extents use many smaller extents instead of a few larger ones. By using this parameter, the default indirect extent size can be increased so large that files in indirects use fewer larger extents. The tunable <code>default_indir_size</code> should be used carefully. If it is set too large, then writes fail when they are unable to allocate extents of the indirect extent size to a file. In general, the fewer and the larger the files on a file system, the larger the <code>default_indir_size</code> can be set. This parameter should generally be set to some multiple of the <code>read_pref_io</code> parameter.</p> <p><code>default_indir_size</code> is not applicable on Version 4 disk layouts.</p>
qio_cache_enable	<p>Enables or disables caching on Quick I/O files. The default behavior is to disable caching. To enable caching, set <code>qio_cache_enable</code> to 1. On systems with large memories, the database cannot always use all of the memory as a cache. By enabling file system caching as a second level cache, performance may be improved. If the database is performing sequential scans of tables, the scans may run faster by enabling file system caching so the file system performs aggressive read-ahead on the files.</p>

Table 2-1 Tunable VxFS I/O parameters (*continued*)

Parameter	Description
<code>read_ahead</code>	<p>The default for all VxFS read operations is to perform sequential read ahead. You can specify the <code>read_ahead</code> cache advisory to implement the VxFS enhanced read ahead functionality. This allows read aheads to detect more elaborate patterns, such as increasing or decreasing read offsets or multithreaded file accesses, in addition to simple sequential reads. You can specify the following values for <code>read_ahead</code>:</p> <ul style="list-style-type: none"> 0—Disables read ahead functionality 1—Retains traditional sequential read ahead behavior 2—Enables enhanced read ahead for all reads <p>The default is 1—VxFS detects only sequential patterns. <code>read_ahead</code> detects patterns on a per-thread basis, up to a maximum determined by <code>vx_era_nthreads</code> parameter. The default number of threads is 5, but you can change the default value by setting the <code>vx_era_nthreads</code> parameter in the <code>/etc/system</code> configuration file.</p>

Table 2-1 Tunable VxFS I/O parameters (continued)

Parameter	Description
<code>write_throttle</code>	<p>The <code>write_throttle</code> parameter is useful in special situations where a computer system has a combination of a large amount of memory and slow storage devices. In this configuration, sync operations, such as <code>fsync()</code>, may take long enough to complete that a system appears to hang. This behavior occurs because the file system is creating dirty buffers (in-memory updates) faster than they can be asynchronously flushed to disk without slowing system performance.</p> <p>Lowering the value of <code>write_throttle</code> limits the number of dirty buffers per file that a file system generates before flushing the buffers to disk. After the number of dirty buffers for a file reaches the <code>write_throttle</code> threshold, the file system starts flushing buffers to disk even if free memory is still available.</p> <p>The default value of <code>write_throttle</code> is zero, which puts no limit on the number of dirty buffers per file. If non-zero, VxFS limits the number of dirty buffers per file to <code>write_throttle</code> buffers.</p> <p>The default value typically generates a large number of dirty buffers, but maintains fast user writes. Depending on the speed of the storage device, if you lower <code>write_throttle</code>, user write performance may suffer, but the number of dirty buffers is limited, so sync operations complete much faster.</p> <p>Because lowering <code>write_throttle</code> may in some cases delay write requests (for example, lowering <code>write_throttle</code> may increase the file disk queue to the <code>max_diskq</code> value, delaying user writes until the disk queue decreases), it is advisable not to change the value of <code>write_throttle</code> unless your system has a combination of large physical memory and slow storage devices.</p>

File system tuning guidelines

If the file system is being used with VxVM, it is advisable to let the VxFS I/O parameters be set to default values based on the volume geometry.

Note: VxFS does not query VxVM with multiple volume sets. To improve I/O performance when using multiple volume sets, use the `vxtunefs` command.

If the file system is being used with a hardware disk array or volume manager other than VxVM, try to align the parameters to match the geometry of the logical disk. With striping or RAID-5, it is common to set `read_pref_io` to the stripe unit size and `read_nstream` to the number of columns in the stripe. For striped arrays, use the same values for `write_pref_io` and `write_nstream`, but for RAID-5 arrays, set `write_pref_io` to the full stripe size and `write_nstream` to 1.

For an application to do efficient disk I/O, it should use the following formula to issue read requests:

■ `read requests = read_nstream x by read_pref_io`

Generally, any multiple or factor of `read_nstream` multiplied by `read_pref_io` should be a good size for performance. For writing, the same rule of thumb applies to the `write_pref_io` and `write_nstream` parameters. When tuning a file system, the best thing to do is try out the tuning parameters under a real life workload.

If an application is performing sequential I/O to large files, the application should try to issue requests larger than `discovered_direct_iosz`. This causes the I/O requests to be performed as discovered direct I/O requests, which are unbuffered like direct I/O but do not require synchronous inode updates when extending the file. If the file is larger than can fit in the cache, using unbuffered I/O avoids removing useful data out of the cache and lessens CPU overhead.

Extent attributes

This chapter includes the following topics:

- [About extent attributes](#)
- [Commands related to extent attributes](#)

About extent attributes

Veritas File System (VxFS) allocates disk space to files in groups of one or more adjacent blocks called extents. VxFS defines an application interface that allows programs to control various aspects of the extent allocation for a given file. The extent allocation policies associated with a file are referred to as extent attributes.

The VxFS `gettext` and `settext` commands let you view or manipulate file extent attributes.

The two basic extent attributes associated with a file are its reservation and its fixed extent size. You can preallocate space to the file by manipulating a file's reservation, or override the default allocation policy of the file system by setting a fixed extent size.

Other policies determine the way these attributes are expressed during the allocation process.

You can specify the following criteria:

- The space reserved for a file must be contiguous
- No allocations will be made for a file beyond the current reservation
- An unused reservation will be released when the file is closed
- Space will be allocated, but no reservation will be assigned
- The file size will be changed to incorporate the allocated space immediately

Some of the extent attributes are persistent and become part of the on-disk information about the file, while other attributes are temporary and are lost after the file is closed or the system is rebooted. The persistent attributes are similar to the file's permissions and are written in the inode for the file. When a file is copied, moved, or archived, only the persistent attributes of the source file are preserved in the new file.

See [“Other controls”](#) on page 65.

In general, the user will only set extent attributes for reservation. Many of the attributes are designed for applications that are tuned to a particular pattern of I/O or disk alignment.

See the `mkfs_vxfs(1M)` manual page.

See [“About VxFS I/O”](#) on page 69.

Reservation: preallocating space to a file

VxFS makes it possible to preallocate space to a file at the time of the request rather than when data is written into the file. This space cannot be allocated to other files in the file system. VxFS prevents any unexpected out-of-space condition on the file system by ensuring that a file's required space will be associated with the file before it is required.

A persistent reservation is not released when a file is truncated. The reservation must be cleared or the file must be removed to free the reserved space.

Fixed extent size

The VxFS default allocation policy uses a variety of methods to determine how to make an allocation to a file when a write requires additional space. The policy attempts to balance the two goals of optimum I/O performance through large allocations and minimal file system fragmentation. VxFS accomplishes these goals by allocating from space available in the file system that best fits the data.

Setting a fixed extent size overrides the default allocation policies for a file and always serves as a persistent attribute. Be careful to choose an extent size appropriate to the application when using fixed extents. An advantage of the VxFS extent-based allocation policies is that they rarely use indirect blocks compared to block based file systems; VxFS eliminates many instances of disk access that stem from indirect references. However, a small extent size can eliminate this advantage.

Files with large extents tend to be more contiguous and have better I/O characteristics. However, the overall performance of the file system degrades because the unused space fragments free space by breaking large extents into

smaller pieces. By erring on the side of minimizing fragmentation for the file system, files may become so non-contiguous that their I/O characteristics would degrade.

Fixed extent sizes are particularly appropriate in the following situations:

- If a file is large and contiguous, a large fixed extent size can minimize the number of extents in the file.

Custom applications may also use fixed extent sizes for specific reasons, such as the need to align extents to cylinder or striping boundaries on disk.

Other controls

The auxiliary controls on extent attributes determine the following conditions:

- Whether allocations are aligned
- Whether allocations are contiguous
- Whether the file can be written beyond its reservation
- Whether an unused reservation is released when the file is closed
- Whether the reservation is a persistent attribute of the file
- When the space reserved for a file will actually become part of the file

Alignment

Specific alignment restrictions coordinate a file's allocations with a particular I/O pattern or disk alignment. Alignment can only be specified if a fixed extent size has also been set. Setting alignment restrictions on allocations is best left to well-designed applications.

See the `mkfs_vxfs(1M)` manual page.

See [“About VxFS I/O”](#) on page 69.

Contiguity

A reservation request can specify that its allocation remain contiguous (all one extent). Maximum contiguity of a file optimizes its I/O characteristics.

Note: Fixed extent sizes or alignment cause a file system to return an error message reporting insufficient space if no suitably sized (or aligned) extent is available. This can happen even if the file system has sufficient free space and the fixed extent size is large.

Write operations beyond reservation

A reservation request can specify that no allocations can take place after a write operation fills the last available block in the reservation. This request can be used a way similar to the function of the `ulimit` command to prevent a file's uncontrolled growth.

Reservation trimming

A reservation request can specify that any unused reservation be released when the file is closed. The file is not completely closed until all processes open against the file have closed it.

Reservation persistence

A reservation request can ensure that the reservation does not become a persistent attribute of the file. The unused reservation is discarded when the file is closed.

Including reservation in the file

A reservation request can make sure the size of the file is adjusted to include the reservation. Normally, the space of the reservation is not included in the file until an extending write operation requires it. A reservation that immediately changes the file size can generate large temporary files. Unlike a `ftruncate` operation that increases the size of a file, this type of reservation does not perform zeroing of the blocks included in the file and limits this facility to users with appropriate privileges. The data that appears in the file may have been previously contained in another file. For users who do not have the appropriate privileges, there is a variant request that prevents such users from viewing uninitialized data.

Commands related to extent attributes

The VxFS commands for manipulating extent attributes are `setext` and `getext`; they allow the user to set up files with a given set of extent attributes or view any attributes that are already associated with a file.

See the `setext(1)` and `getext(1M)` manual pages.

The VxFS-specific commands `vxdump` and `vxrestore` preserve extent attributes when backing up, restoring, moving, or copying files.

Most of these commands include a command line option (`-e`) for maintaining extent attributes on files. This option specifies dealing with a VxFS file that has extent attribute information including reserved space, a fixed extent size, and extent alignment. The extent attribute information may be lost if the destination

file system does not support extent attributes, has a different block size than the source file system, or lacks free extents appropriate to satisfy the extent attribute requirements.

The `-e` option takes any of the following keywords as an argument:

<code>warn</code>	Issues a warning message if extent attribute information cannot be maintained (the default)
<code>force</code>	Fails the copy if extent attribute information cannot be maintained
<code>ignore</code>	Ignores extent attribute information entirely

Failure to preserve extent attributes

Whenever a file is copied, moved, or archived using commands that preserve extent attributes, there is nevertheless the possibility of losing the attributes.

Such a failure might occur for one of the following reasons:

- The file system receiving a copied, moved, or restored file from an archive is not a VxFS type. Since other file system types do not support the extent attributes of the VxFS file system, the attributes of the source file are lost during the migration.
- The file system receiving a copied, moved, or restored file is a VxFS type but does not have enough free space to satisfy the extent attributes. For example, consider a 50K file and a reservation of 1 MB. If the target file system has 500K free, it could easily hold the file but fail to satisfy the reservation.
- The file system receiving a copied, moved, or restored file from an archive is a VxFS type but the different block sizes of the source and target file system make extent attributes impossible to maintain. For example, consider a source file system of block size 1024, a target file system of block size 4096, and a file that has a fixed extent size of 3 blocks (3072 bytes). This fixed extent size adapts to the source file system but cannot translate onto the target file system. The same source and target file systems in the preceding example with a file carrying a fixed extent size of 4 could preserve the attribute; a 4 block (4096 byte) extent on the source file system would translate into a 1 block extent on the target.

On a system with mixed block sizes, a copy, move, or restoration operation may or may not succeed in preserving attributes. It is recommended that the same block size be used for all file systems on a given system.

VxFS I/O Overview

This chapter includes the following topics:

- [About VxFS I/O](#)
- [Buffered and Direct I/O](#)
- [Concurrent I/O](#)
- [Cache advisories](#)
- [Freezing and thawing a file system](#)
- [Getting the I/O size](#)

About VxFS I/O

VxFS processes two basic types of file system I/O:

- Sequential
- Random or I/O that is not sequential

For sequential I/O, VxFS employs a read-ahead policy by default when the application is reading data. For writing, it allocates contiguous blocks if possible. In most cases, VxFS handles I/O that is sequential through buffered I/O. VxFS handles random or nonsequential I/O using direct I/O without buffering.

VxFS provides a set of I/O cache advisories for use when accessing files.

See the *Veritas File System Programmer's Reference Guide*.

See the `vxfsio(7)` manual page.

Buffered and Direct I/O

VxFS responds with read-ahead for sequential read I/O. This results in buffered I/O. The data is prefetched and retained in buffers for the application. The data buffers are commonly referred to as VxFS buffer cache. This is the default VxFS behavior.

On the other hand, direct I/O does not buffer the data when the I/O to the underlying device is completed. This saves system resources like memory and CPU usage. Direct I/O is possible only when alignment and sizing criteria are satisfied.

See “[Direct I/O requirements](#)” on page 70.

All of the supported platforms have a VxFS buffered cache. Each platform also has either a page cache or its own buffer cache. These caches are commonly known as the file system caches.

Direct I/O does not use these caches. The memory used for direct I/O is discarded after the I/O is complete,

Direct I/O

Direct I/O is an unbuffered form of I/O. If the `VX_DIRECT` advisory is set, the user is requesting direct data transfer between the disk and the user-supplied buffer for reads and writes. This bypasses the kernel buffering of data, and reduces the CPU overhead associated with I/O by eliminating the data copy between the kernel buffer and the user's buffer. This also avoids taking up space in the buffer cache that might be better used for something else. The direct I/O feature can provide significant performance gains for some applications.

The direct I/O and `VX_DIRECT` advisories are maintained on a per-file-descriptor basis.

Direct I/O requirements

For an I/O operation to be performed as direct I/O, it must meet certain alignment criteria. The alignment constraints are usually determined by the disk driver, the disk controller, and the system memory management hardware and software.

The requirements for direct I/O are as follows:

- The starting file offset must be aligned to a 512-byte boundary.
- The ending file offset must be aligned to a 512-byte boundary, or the length must be a multiple of 512 bytes.
- The memory buffer must start on an 8-byte boundary.

Direct I/O versus synchronous I/O

Because direct I/O maintains the same data integrity as synchronous I/O, it can be used in many applications that currently use synchronous I/O. If a direct I/O request does not allocate storage or extend the file, the inode is not immediately written.

Direct I/O CPU overhead

The CPU cost of direct I/O is about the same as a raw disk transfer. For sequential I/O to very large files, using direct I/O with large transfer sizes can provide the same speed as buffered I/O with much less CPU overhead.

If the file is being extended or storage is being allocated, direct I/O must write the inode change before returning to the application. This eliminates some of the performance advantages of direct I/O.

Discovered Direct I/O

Discovered Direct I/O is a file system tunable that is set using the `vxtunefs` command. When the file system gets an I/O request larger than the `discovered_direct_iosz`, it tries to use direct I/O on the request. For large I/O sizes, Discovered Direct I/O can perform much better than buffered I/O.

Discovered Direct I/O behavior is similar to direct I/O and has the same alignment constraints, except writes that allocate storage or extend the file size do not require writing the inode changes before returning to the application.

See [“Tuning I/O”](#) on page 51.

Unbuffered I/O

If the `VX_UNBUFFERED` advisory is set, I/O behavior is the same as direct I/O with the `VX_DIRECT` advisory set, so the alignment constraints that apply to direct I/O also apply to unbuffered I/O. For unbuffered I/O, however, if the file is being extended, or storage is being allocated to the file, inode changes are not updated synchronously before the write returns to the user. The `VX_UNBUFFERED` advisory is maintained on a per-file-descriptor basis.

See [“Tuning I/O”](#) on page 51.

Data synchronous I/O

If the `VX_DSYNC` advisory is set, the user is requesting data synchronous I/O. In synchronous I/O, the data is written, and the inode is written with updated times and, if necessary, an increased file size. In data synchronous I/O, the data is

transferred to disk synchronously before the write returns to the user. If the file is not extended by the write, the times are updated in memory, and the call returns to the user. If the file is extended by the operation, the inode is written before the write returns.

The direct I/O and `VX_DSYNC` advisories are maintained on a per-file-descriptor basis.

Data synchronous I/O vs. synchronous I/O

Like direct I/O, the data synchronous I/O feature can provide significant application performance gains. Because data synchronous I/O maintains the same data integrity as synchronous I/O, it can be used in many applications that currently use synchronous I/O. If the data synchronous I/O does not allocate storage or extend the file, the inode is not immediately written. The data synchronous I/O does not have any alignment constraints, so applications that find it difficult to meet the alignment constraints of direct I/O should use data synchronous I/O.

If the file is being extended or storage is allocated, data synchronous I/O must write the inode change before returning to the application. This case eliminates the performance advantage of data synchronous I/O.

Concurrent I/O

Concurrent I/O (`VX_CONCURRENT`) allows multiple processes to read from or write to the same file without blocking other `read(2)` or `write(2)` calls. POSIX semantics requires `read` and `write` calls to be serialized on a file with other `read` and `write` calls. With POSIX semantics, a `read` call either reads the data before or after the `write` call occurred. With the `VX_CONCURRENT` advisory set, the `read` and `write` operations are not serialized as in the case of a character device. This advisory is generally used by applications that require high performance for accessing data and do not perform overlapping writes to the same file. It is the responsibility of the application or the running threads to coordinate the `write` activities to the same file when using Concurrent I/O.

Concurrent I/O can be enabled in the following ways:

- By specifying the `VX_CONCURRENT` advisory flag for the file descriptor in the `VX_SETCACHE` `ioctl` command. Only the `read(2)` and `write(2)` calls occurring through this file descriptor use concurrent I/O. The `read` and `write` operations occurring through other file descriptors for the same file will still follow the POSIX semantics.

See `vxfsio(7)` manual page.

- By using the `cio` mount option. The `read(2)` and `write(2)` operations occurring on all of the files in this particular file system will use concurrent I/O. See “[The cio option](#)” on page 43. See the `mount_vxfs(1M)` manual page.

Cache advisories

VxFS allows an application to set cache advisories for use when accessing files. VxFS cache advisories enable applications to help monitor the buffer cache and provide information on how better to tune the buffer cache to improve performance gain.

The basic function of the cache advisory is to let you know whether you could have avoided a later re-read of block X if the buffer cache had been a little larger. Conversely, the cache advisory can also let you know that you could safely reduce the buffer cache size without putting block X into jeopardy.

These advisories are in memory only and do not persist across reboots. Some advisories are currently maintained on a per-file, not a per-file-descriptor, basis. Only one set of advisories can be in effect for all accesses to the file. If two conflicting applications set different advisories, both must use the advisories that were last set.

All advisories are set using the `VX_SETCACHE` ioctl command. The current set of advisories can be obtained with the `VX_GETCACHE` ioctl command.

See the `vxfsio(7)` manual page.

Freezing and thawing a file system

Freezing a file system is a necessary step for obtaining a stable and consistent image of the file system at the volume level. Consistent volume-level file system images can be obtained and used with a file system snapshot tool. The freeze operation flushes all buffers and pages in the file system cache that contain dirty metadata and user data. The operation then suspends any new activity on the file system until the file system is thawed.

The `VX_FREEZE` ioctl command is used to freeze a file system. Freezing a file system temporarily blocks all I/O operations to a file system and then performs a sync on the file system. When the `VX_FREEZE` ioctl is issued, all access to the file system is blocked at the system call level. Current operations are completed and the file system is synchronized to disk.

When the file system is frozen, any attempt to use the frozen file system, except for a `VX_THAW` ioctl command, is blocked until a process executes the `VX_THAW` ioctl command or the time-out on the freeze expires.

Getting the I/O size

VxFS provides the `VX_GET_IOPARAMETERS` ioctl to get the recommended I/O sizes to use on a file system. This ioctl can be used by the application to make decisions about the I/O sizes issued to VxFS for a file or file device.

See the `vxtunefs(1M)` and `vxfsio(7)` manual pages.

See [“Tuning I/O”](#) on page 51.

Online backup using file system snapshots

This chapter includes the following topics:

- [About snapshot file systems](#)
- [Snapshot file system backups](#)
- [Creating a snapshot file system](#)
- [Backup examples](#)
- [Snapshot file system performance](#)
- [Differences Between Snapshots and Storage Checkpoints](#)
- [About snapshot file system disk structure](#)
- [How a snapshot file system works](#)

About snapshot file systems

A snapshot file system is an exact image of a VxFS file system, referred to as the snapped file system, that provides a mechanism for making backups. The snapshot is a consistent view of the file system “snapped” at the point in time the snapshot is made. You can select files to back up from the snapshot using a standard utility such as `cpio` or `cp`, or back up the entire file system image using the `vxdump` or `fscat` utilities.

You use the `mount` command to create a snapshot file system; the `mkfs` command is not required. A snapshot file system is always read-only. A snapshot file system exists only as long as it and the snapped file system are mounted and ceases to exist when unmounted. A snapped file system cannot be unmounted until all of

its snapshots are unmounted. Although it is possible to have multiple snapshots of a file system made at different times, it is not possible to make a snapshot of a snapshot.

Note: A snapshot file system ceases to exist when unmounted. If mounted again, it is actually a fresh snapshot of the snapped file system. A snapshot file system must be unmounted before its dependent snapped file system can be unmounted. Neither the `fuser` command nor the `mount` command will indicate that a snapped file system cannot be unmounted because a snapshot of it exists.

On cluster file systems, snapshots can be created on any node in the cluster, and backup operations can be performed from that node. The snapshot of a cluster file system is accessible only on the node where it is created, that is, the snapshot file system itself cannot be cluster mounted.

See the *Veritas Storage Foundation Cluster File System Administrator's Guide*.

Snapshot file system backups

After a snapshot file system is created, the snapshot maintains a consistent backup of data in the snapped file system.

Backup programs, such as `cpio`, that back up a standard file system tree can be used without modification on a snapshot file system because the snapshot presents the same data as the snapped file system. Backup programs, such as `vxdump`, that access the disk structures of a file system require some modifications to handle a snapshot file system.

VxFS utilities recognize snapshot file systems and modify their behavior so that they operate the same way on snapshots as they do on standard file systems. Other backup programs that typically read the raw disk image cannot work on snapshots without altering the backup procedure.

These other backup programs can use the `fscat` command to obtain a raw image of the entire file system that is identical to an image obtainable by running a `dd` command on the disk device containing the snapped file system at the exact moment the snapshot was created. The `snapread ioctl` takes arguments similar to those of the `read` system call and returns the same results that are obtainable by performing a read on the disk device containing the snapped file system at the exact time the snapshot was created. In both cases, however, the snapshot file system provides a consistent image of the snapped file system with all activity complete—it is an instantaneous read of the entire file system. This is much different than the results that would be obtained by a `dd` or `read` command on the disk device of an active file system.

Creating a snapshot file system

You create a snapshot file system by using the `-o snapof=` option of the `mount` command. The `-o snapsize=` option may also be required if the device you are mounting does not identify the device size in its disk label, or if you want a size smaller than the entire device.

You must make the snapshot file system large enough to hold any blocks on the snapped file system that may be written to while the snapshot file system exists. If a snapshot runs out of blocks to hold copied data, the snapshot is disabled and further attempts to access the snapshot file system fail.

During periods of low activity (such as nights and weekends), a snapshot typically requires about two to six percent of the blocks of the snapped file system. During a period of high activity, the snapshot of a typical file system may require 15 percent of the blocks of the snapped file system. Most file systems do not turn over 15 percent of data in a single day. These approximate percentages tend to be lower for larger file systems and higher for smaller file systems. You can allocate blocks to a snapshot based on characteristics such as file system usage and duration of backups.

Warning: Any existing data on the device used for the snapshot is overwritten.

To create a snapshot file system

- ◆ Mount the file system with the `-o snapof=` option:

```
# mount -F vxfs -o snapof=special,snapsize=snapshot_size \  
snapshot_special snapshot_mount_point
```

Backup examples

In the following examples, the `vxdump` utility is used to ascertain whether `/dev/vx/dsk/fsvol/vol1` is a snapshot mounted as `/backup/home` and does the appropriate work to get the snapshot data through the mount point.

These are typical examples of making a backup of a 300,000 block file system named `/home` using a snapshot file system on `/dev/vx/dsk/fsvol/vol1` with a snapshot mount point of `/backup/home`.

To create a backup using a snapshot file system

- 1 To back up files changed within the last week using `cpio`:

```
# mount -F vxfs -o snapof=/home,snapsize=100000 \
/dev/vx/dsk/fsvol/vol1 /backup/home
# cd /backup
# find home -ctime -7 -depth -print | cpio -oc > /dev/rmt/0m
# umount /backup/home
```

- 2 To do a level 3 backup of `/dev/vx/dsk/fsvol/vol1` and collect those files that have changed in the current directory:

```
# vxdump 3f - /dev/vx/dsk/fsvol/vol1 | vxrestore -xf -
```

- 3 To do a full backup of `/home`, which exists on disk `/dev/vx/dsk/fsvol/vol1`, and use `dd` to control blocking of output onto tape device using `vxdump`:

```
# mount -F vxfs -o snapof=/home,snapsize=100000 \
/dev/vx/dsk/fsvol/vol1 /backup/home
# vxdump f - /dev/vx/dsk/fsvol/vol1 | dd bs=128k > \
/dev/rmt/0m
```

Snapshot file system performance

Snapshot file systems maximize the performance of the snapshot at the expense of writes to the snapped file system. Reads from a snapshot file system typically perform at nearly the throughput rates of reads from a standard VxFS file system.

The performance of reads from the snapped file system are generally not affected. However, writes to the snapped file system, typically average two to three times as long as without a snapshot. This is because the initial write to a data block requires reading the old data, writing the data to the snapshot, and then writing the new data to the snapped file system. If there are multiple snapshots of the same snapped file system, writes are even slower. Only the initial write to a block experiences this delay, so operations such as writes to the intent log or inode updates proceed at normal speed after the initial write.

Reads from the snapshot file system are impacted if the snapped file system is busy because the snapshot reads are slowed by the disk I/O associated with the snapped file system.

The overall impact of the snapshot is dependent on the read to write ratio of an application and the mixing of the I/O operations. For example, a database

application running an online transaction processing (OLTP) workload on a snapped file system was measured at about 15 to 20 percent slower than a file system that was not snapped.

Differences Between Snapshots and Storage Checkpoints

While snapshots and Storage Checkpoints both create a point-in-time image of a file system and only the changed data blocks are updated, there are significant differences between the two technologies:

Table 5-1 Differences between snapshots and Storage Checkpoints

Snapshots	Storage Checkpoints
Require a separate device for storage	Reside on the same device as the original file system
Are read-only	Can be read-only or read-write
Are transient	Are persistent
Cease to exist after being unmounted	Can exist and be mounted on their own
Track changed blocks on the file system level	Track changed blocks on each file in the file system

Storage Checkpoints also serve as the enabling technology for two other Veritas features: Block-Level Incremental Backups and Storage Rollback, which are used extensively for backing up databases.

See [“About Storage Checkpoints”](#) on page 83.

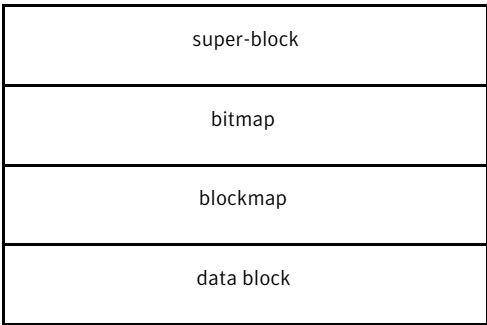
About snapshot file system disk structure

A snapshot file system consists of:

- A super-block
- A bitmap
- A blockmap
- Data blocks copied from the snapped file system

The following figure shows the disk structure of a snapshot file system.

Figure 5-1 The Snapshot Disk Structure



The super-block is similar to the super-block of a standard VxFS file system, but the magic number is different and many of the fields are not applicable.

The bitmap contains one bit for every block on the snapped file system. Initially, all bitmap entries are zero. A set bit indicates that the appropriate block was copied from the snapped file system to the snapshot. In this case, the appropriate position in the blockmap references the copied block.

The blockmap contains one entry for each block on the snapped file system. Initially, all entries are zero. When a block is copied from the snapped file system to the snapshot, the appropriate entry in the blockmap is changed to contain the block number on the snapshot file system that holds the data from the snapped file system.

The data blocks are filled by data copied from the snapped file system, starting from the beginning of the data block area.

How a snapshot file system works

A snapshot file system is created by mounting an empty disk slice as a snapshot of a currently mounted file system. The bitmap, blockmap and super-block are initialized and then the currently mounted file system is frozen. After the file system to be snapped is frozen, the snapshot is enabled and mounted and the snapped file system is thawed. The snapshot appears as an exact image of the snapped file system at the time the snapshot was made.

See [“Freezing and thawing a file system”](#) on page 73.

Initially, the snapshot file system satisfies read requests by finding the data on the snapped file system and returning it to the requesting process. When an inode update or a write changes the data in block n of the snapped file system, the old data is first read and copied to the snapshot before the snapped file system is updated. The bitmap entry for block n is changed from 0 to 1, indicating that the

data for block *n* can be found on the snapshot file system. The blockmap entry for block *n* is changed from 0 to the block number on the snapshot file system containing the old data.

A subsequent read request for block *n* on the snapshot file system will be satisfied by checking the bitmap entry for block *n* and reading the data from the indicated block on the snapshot file system, instead of from block *n* on the snapped file system. This technique is called copy-on-write. Subsequent writes to block *n* on the snapped file system do not result in additional copies to the snapshot file system, since the old data only needs to be saved once.

All updates to the snapped file system for inodes, directories, data in files, extent maps, and so forth, are handled in this fashion so that the snapshot can present a consistent view of all file system structures on the snapped file system for the time when the snapshot was created. As data blocks are changed on the snapped file system, the snapshot gradually fills with data copied from the snapped file system.

The amount of disk space required for the snapshot depends on the rate of change of the snapped file system and the amount of time the snapshot is maintained. In the worst case, the snapped file system is completely full and every file is removed and rewritten. The snapshot file system would need enough blocks to hold a copy of every block on the snapped file system, plus additional blocks for the data structures that make up the snapshot file system. This is approximately 101 percent of the size of the snapped file system. Normally, most file systems do not undergo changes at this extreme rate. During periods of low activity, the snapshot should only require two to six percent of the blocks of the snapped file system. During periods of high activity, the snapshot might require 15 percent of the blocks of the snapped file system. These percentages tend to be lower for larger file systems and higher for smaller ones.

Warning: If a snapshot file system runs out of space for changed data blocks, it is disabled and all further attempts to access it fails. This does not affect the snapped file system.

Storage Checkpoints

This chapter includes the following topics:

- [About Storage Checkpoints](#)
- [How a Storage Checkpoint works](#)
- [Types of Storage Checkpoints](#)
- [Storage Checkpoint administration](#)
- [Space management considerations](#)
- [Restoring a file system from a Storage Checkpoint](#)
- [Storage Checkpoint quotas](#)

About Storage Checkpoints

Veritas File System provides a Storage Checkpoint feature that quickly creates a persistent image of a file system at an exact point in time. Storage Checkpoints significantly reduce I/O overhead by identifying and maintaining only the file system blocks that have changed since the last Storage Checkpoint or backup via a copy-on-write technique.

See “[Copy-on-write](#)” on page 87.

Storage Checkpoints provide:

- Persistence through reboots and crashes.
- The ability for data to be immediately writeable by preserving the file system metadata, the directory hierarchy, and user data.

Storage Checkpoints are actually data objects that are managed and controlled by the file system. You can create, remove, and rename Storage Checkpoints because they are data objects with associated names.

See [“How a Storage Checkpoint works”](#) on page 85.

Unlike a disk-based mirroring technology that requires a separate storage space, Storage Checkpoints minimize the use of disk space by using a Storage Checkpoint within the same free space available to the file system.

After you create a Storage Checkpoint of a mounted file system, you can also continue to create, remove, and update files on the file system without affecting the logical image of the Storage Checkpoint. A Storage Checkpoint preserves not only the name space (directory hierarchy) of the file system, but also the user data as it existed at the moment the file system image was captured.

You can use a Storage checkpoint in many ways. For example, you can use them to:

- Create a stable image of the file system that can be backed up to tape.
- Provide a mounted, on-disk backup of the file system so that end users can restore their own files in the event of accidental deletion. This is especially useful in a home directory, engineering, or email environment.
- Create a copy of an application's binaries before installing a patch to allow for rollback in case of problems.
- Create an on-disk backup of the file system in that can be used addition to a traditional tape-based backup to provide faster backup and restore capabilities.

How Storage Checkpoints differ from snapshots

Storage Checkpoints differ from Veritas File System snapshots in the following ways because they:

- Allow write operations to the Storage Checkpoint itself.
- Persist after a system reboot or failure.
- Share the same pool of free space as the file system.
- Maintain a relationship with other Storage Checkpoints by identifying changed file blocks since the last Storage Checkpoint.
- Have multiple, read-only Storage Checkpoints that reduce I/O operations and required storage space because the most recent Storage Checkpoint is the only one that accumulates updates from the primary file system.

Various backup and replication solutions can take advantage of Storage Checkpoints. The ability of Storage Checkpoints to track the file system blocks that have changed since the last Storage Checkpoint facilitates backup and replication applications that only need to retrieve the changed data. Storage Checkpoints significantly minimize data movement and may promote higher

availability and data integrity by increasing the frequency of backup and replication solutions.

Storage Checkpoints can be taken in environments with a large number of files, such as file servers with millions of files, with little adverse impact on performance. Because the file system does not remain frozen during Storage Checkpoint creation, applications can access the file system even while the Storage Checkpoint is taken. However, Storage Checkpoint creation may take several minutes to complete depending on the number of files in the file system.

How a Storage Checkpoint works

The Storage Checkpoint facility freezes the mounted file system (known as the primary fileset), initializes the Storage Checkpoint, and thaws the file system. Specifically, the file system is first brought to a stable state where all of its data is written to disk, and the freezing process momentarily blocks all I/O operations to the file system. A Storage Checkpoint is then created without any actual data; the Storage Checkpoint instead points to the block map of the primary fileset. The thawing process that follows restarts I/O operations to the file system.

You can create a Storage Checkpoint on a single file system or a list of file systems. A Storage Checkpoint of multiple file systems simultaneously freezes the file systems, creates a Storage Checkpoint on all of the file systems, and thaws the file systems. As a result, the Storage Checkpoints for multiple file systems have the same creation timestamp. The Storage Checkpoint facility guarantees that multiple file system Storage Checkpoints are created on all or none of the specified file systems, unless there is a system crash while the operation is in progress.

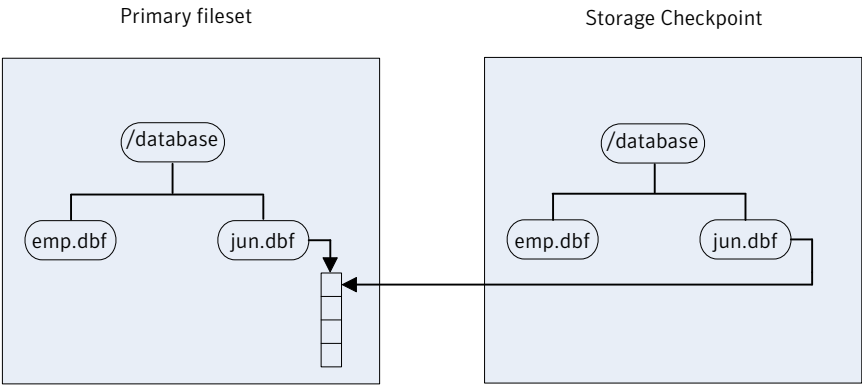
Note: The calling application is responsible for cleaning up Storage Checkpoints after a system crash.

A Storage Checkpoint of the primary fileset initially contains a pointer to the file system block map rather than to any actual data. The block map points to the data on the primary fileset.

Figure 6-1 shows the file system `/database` and its Storage Checkpoint.

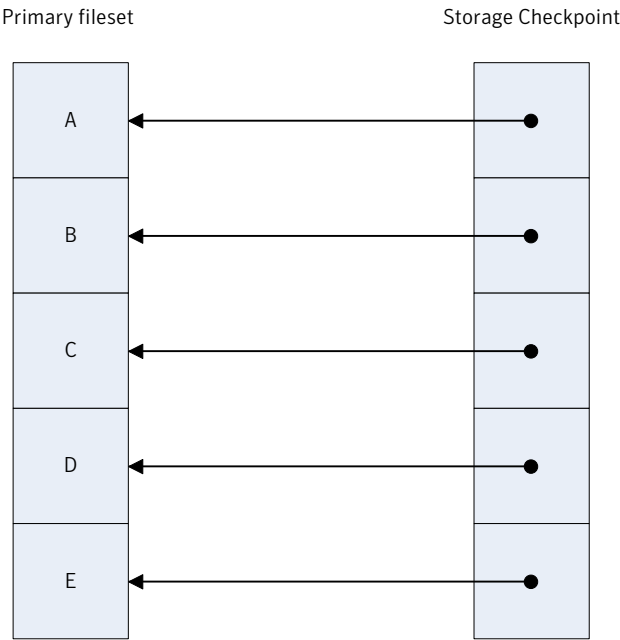
The Storage Checkpoint is logically identical to the primary fileset when the Storage Checkpoint is created, but it does not contain any actual data blocks.

Figure 6-1 Primary fileset and its Storage Checkpoint



In [Figure 6-2](#), a square represents each block of the file system. This figure shows a Storage Checkpoint containing pointers to the primary fileset at the time the Storage Checkpoint is taken, as in [Figure 6-1](#).

Figure 6-2 Initializing a Storage Checkpoint



The Storage Checkpoint presents the exact image of the file system by finding the data from the primary fileset. As the primary fileset is updated, the original

data is copied to the Storage Checkpoint before the new data is written. When a write operation changes a specific data block in the primary fileset, the old data is first read and copied to the Storage Checkpoint before the primary fileset is updated. Subsequent writes to the specified data block on the primary fileset do not result in additional updates to the Storage Checkpoint because the old data needs to be saved only once. As blocks in the primary fileset continue to change, the Storage Checkpoint accumulates the original data blocks.

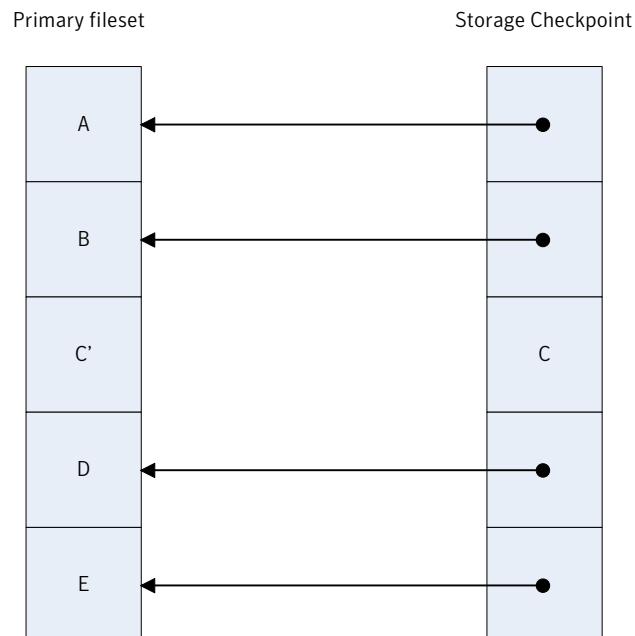
Copy-on-write

In [Figure 6-3](#), the third block originally containing C is updated.

Before the block is updated with new data, the original data is copied to the Storage Checkpoint. This is called the copy-on-write technique, which allows the Storage Checkpoint to preserve the image of the primary fileset when the Storage Checkpoint is taken.

Every update or write operation does not necessarily result in the process of copying data to the Storage Checkpoint. In this example, subsequent updates to this block, now containing C', are not copied to the Storage Checkpoint because the original image of the block containing C is already saved.

Figure 6-3 Updates to the primary fileset



Types of Storage Checkpoints

You can create the following types of Storage Checkpoints:

- [Data Storage Checkpoints](#)
- [nodata Storage Checkpoints](#)
- [Removable Storage Checkpoints](#)
- [Non-mountable Storage Checkpoints](#)

Data Storage Checkpoints

A data Storage Checkpoint is a complete image of the file system at the time the Storage Checkpoint is created. This type of Storage Checkpoint contains the file system metadata and file data blocks. You can mount, access, and write to a data Storage Checkpoint just as you would to a file system. Data Storage Checkpoints are useful for backup applications that require a consistent and stable image of an active file system. Data Storage Checkpoints introduce some overhead to the system and to the application performing the write operation. For best results, limit the life of data Storage Checkpoints to minimize the impact on system resources.

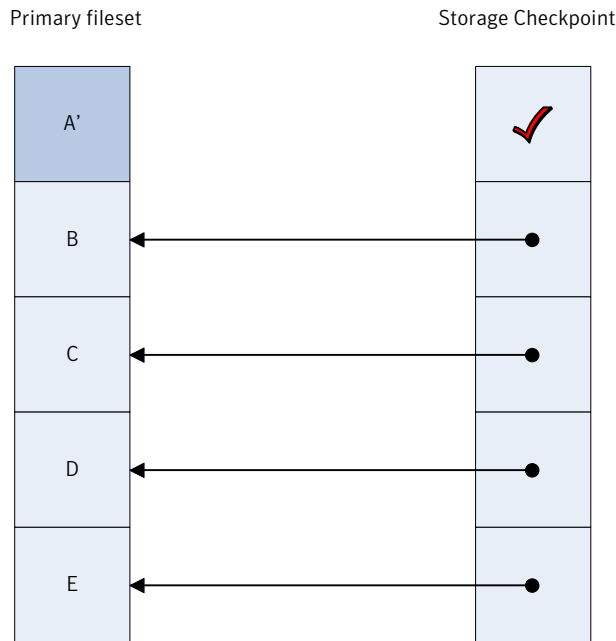
See [“Showing the difference between a data and a nodata Storage Checkpoint”](#) on page 95.

nodata Storage Checkpoints

A nodata Storage Checkpoint only contains file system metadata—no file data blocks. As the original file system changes, the nodata Storage Checkpoint records the location of every changed block. Nodata Storage Checkpoints use minimal system resources and have little impact on the performance of the file system because the data itself does not have to be copied.

In [Figure 6-4](#), the first block originally containing A is updated.

The original data is not copied to the storage checkpoint, but the changed block is marked in the Storage Checkpoint. The marker indicates which data has changed.

Figure 6-4 Updates to a nodata clone

See [“Showing the difference between a data and a nodata Storage Checkpoint”](#) on page 95.

Removable Storage Checkpoints

A removable Storage Checkpoint can “self-destruct” under certain conditions when the file system runs out of space.

See [“Space management considerations”](#) on page 101.

After encountering certain out-of-space (`ENOSPC`) conditions, the kernel removes Storage Checkpoints to free up space for the application to continue running on the file system. In almost all situations, you should create Storage Checkpoints with the removable attribute.

Non-mountable Storage Checkpoints

You can create Storage Checkpoints that cannot be mounted by using the `fsckptadm set nomount` command.

See the `fsckptadm(1M)` manual page.

Use this type of Storage Checkpoint as a security feature which prevents other applications from accessing the Storage Checkpoint and modifying it.

Storage Checkpoint administration

Storage Checkpoint administrative operations require the `fsckptadm` utility.

See the `fsckptadm(1M)` manual page.

You can use the `fsckptadm` utility to create and remove Storage Checkpoints, change attributes, and ascertain statistical data. Every Storage Checkpoint has an associated name, which allows you to manage Storage Checkpoints; this name is limited to 127 characters and cannot contain a colon (:).

Storage Checkpoints require some space for metadata on the volume or set of volumes specified by the file system allocation policy or Storage Checkpoint allocation policy. The `fsckptadm` utility displays an error if the volume or set of volumes does not have enough free space to contain the metadata. You can roughly approximate the amount of space required by the metadata using a method that depends on the disk layout version of the file system.

For disk layout Version 5 or prior, multiply the number of inodes (# of inodes) by the inode size (*inosize*) in bytes, and add 1 or 2 megabytes to get the approximate amount of space required. You can determine the number of inodes with the `fsckptadm` utility, and the inode size with the `mkfs` command:

```
# fsckptadm -v info ' ' /mnt0
UNNAMED:
    ctime      = Thu 3 Mar 2005 7:00:17 PM PST
    mtime      = Thu 3 Mar 2005 7:00:17 PM PST
    flags      = largefiles, mounted
# of inodes    = 23872
# of blocks    = 27867
.
.
.
# of overlay bmaps = 0
# mkfs -m /dev/vx/rdisk/sharedg/vol0
# mkfs -F vxfs -o \
    bsize=1024,version=5,inosize=256,logsize=65536,\
    largefiles /dev/vx/rdisk/sharedg/vol0
```

In this example, the approximate amount of space required by the metadata is 7 or 8 megabytes (23,872 x 256 bytes, plus 1 or 2 megabytes).

For disk layout Version 6 or 7, multiply the number of inodes by 1 byte, and add 1 or 2 megabytes to get the approximate amount of space required. You can determine the number of inodes with the `fsckptadm` utility as above. Using the output from the example for disk layout Version 5, the approximate amount of space required by the metadata is just over one or two megabytes (23,872 x 1 byte, plus 1 or 2 megabytes).

Use the `fsvoladm` command to determine if the volume set has enough free space.

See the `fsvoladm(1M)` manual page.

```
# fsvoladm list /mnt0
devid      size      used      avail      name
0          20971520      8497658      12473862      mnt1
1          20971520      6328993      14642527      mnt2
2          20971520      4458462      16513058      mnt3
```

Creating a Storage Checkpoint

The following example shows the creation of a nodata Storage Checkpoint named `thu_7pm` on `/mnt0` and lists all Storage Checkpoints of the `/mnt0` file system:

```
# fsckptadm -n create thu_7pm /mnt0
# fsckptadm list /mnt0
/mnt0
thu_7pm:
  ctime   = Thu 3 Mar 2005 7:00:17 PM PST
  mtime   = Thu 3 Mar 2005 7:00:17 PM PST
  flags   = nodata, largefiles
```

The following example shows the creation of a removable Storage Checkpoint named `thu_8pm` on `/mnt0` and lists all Storage Checkpoints of the `/mnt0` file system:

```
# fsckptadm -r create thu_8pm /mnt0
# fsckptadm list /mnt0
/mnt0
thu_8pm:
  ctime   = Thu 3 Mar 2005 8:00:19 PM PST
  mtime   = Thu 3 Mar 2005 8:00:19 PM PST
  flags   = largefiles, removable
thu_7pm:
  ctime   = Thu 3 Mar 2005 7:00:17 PM PST
  mtime   = Thu 3 Mar 2005 7:00:17 PM PST
  flags   = nodata, largefiles
```

Removing a Storage Checkpoint

You can delete a Storage Checkpoint by specifying the `remove` keyword of the `fsckptadm` command. Specifically, you can use either the synchronous or asynchronous method of removing a Storage Checkpoint; the asynchronous method is the default method. The synchronous method entirely removes the Storage Checkpoint and returns all of the blocks to the file system before completing the `fsckptadm` operation. The asynchronous method simply marks the Storage Checkpoint for removal and causes `fsckptadm` to return immediately. At a later time, an independent kernel thread completes the removal operation and releases the space used by the Storage Checkpoint.

In this example, `/mnt0` is a mounted VxFS file system with a Version 6 disk layout. This example shows the asynchronous removal of the Storage Checkpoint named `thu_8pm` and synchronous removal of the Storage Checkpoint named `thu_7pm`. This example also lists all the Storage Checkpoints remaining on the `/mnt0` file system after the specified Storage Checkpoint is removed:

```
# fsckptadm remove thu_8pm /mnt0
# fsckptadm list /mnt0
/mnt0
thu_7pm:
  ctime   = Thu 3 Mar 2005 7:00:17 PM PST
  mtime   = Thu 3 Mar 2005 7:00:17 PM PST
  flags   = nodata, largefiles
# fsckptadm -s remove thu_7pm /mnt0
# fsckptadm list /mnt0
/mnt0
```

Accessing a Storage Checkpoint

You can mount Storage Checkpoints using the `mount` command with the `mount` option `-o ckpt=ckpt_name`.

See the `mount_vxfs(1M)` manual page.

Observe the following rules when mounting Storage Checkpoints:

- Storage Checkpoints are mounted as read-only Storage Checkpoints by default. If you must write to a Storage Checkpoint, mount it using the `-o rw` option.
- If a Storage Checkpoint is currently mounted as a read-only Storage Checkpoint, you can remount it as a writable Storage Checkpoint using the `-o remount` option.

- To mount a Storage Checkpoint of a file system, first mount the file system itself.
- To unmount a file system, first unmount all of its Storage Checkpoints.

Warning: If you create a Storage Checkpoint for backup purposes, do not mount it as a writable Storage Checkpoint. You will lose the point-in-time image if you accidentally write to the Storage Checkpoint.

A Storage Checkpoint is mounted on a special pseudo device. This pseudo device does not exist in the system name space; the device is internally created by the system and used while the Storage Checkpoint is mounted. The pseudo device is removed after you unmount the Storage Checkpoint. A pseudo device name is formed by appending the Storage Checkpoint name to the file system device name using the colon character (:) as the separator.

For example, if a Storage Checkpoint named `may_23` belongs to the file system residing on the special device `/dev/vx/dsk/fsvol/vol1`, the Storage Checkpoint pseudo device name is:

```
/dev/vx/dsk/fsvol/vol1:may_23
```

- To mount the Storage Checkpoint named `may_23` as a read-only Storage Checkpoint on directory `/fsvol_may_23`, type:

```
# mount -F vxfs -o ckpt=may_23 /dev/vx/dsk/fsvol/vol1:may_23 \
/fsvol_may_23
```

Note: The `vol1` file system must already be mounted before the Storage Checkpoint can be mounted.

- To remount the Storage Checkpoint named `may_23` as a writable Storage Checkpoint, type:

```
# mount -F vxfs -o ckpt=may_23,remount \
/dev/vx/dsk/fsvol/vol1:may_23 /fsvol_may_23
```

- To mount this Storage Checkpoint automatically when the system starts up, put the following entries in the `/etc/fstab` file:

Device-Special-File	Mount-Point	fstype	options	backup- frequency	pas- numl
/dev/vx/dsk/fsvol/	/fsvol	vxfs	defaults	0	0

```
vol1  
/dev/vx/dsk/fsvol/ /fsvol_may_23 vxfs ckpt=may_23 0  
vol1:may_23
```

0

- To mount a Storage Checkpoint of a cluster file system, you must also use the `-o cluster` option:

```
# mount -F vxfs -o cluster,ckpt=may_23 \  
/dev/vx/dsk/fsvol/vol1:may_23 /fsvol_may_23
```

You can only mount a Storage Checkpoint cluster-wide if the file system that the Storage Checkpoint belongs to is also mounted cluster-wide. Similarly, you can only mount a Storage Checkpoint locally if the file system that the Storage Checkpoint belongs to is mounted locally.

You can unmount Storage Checkpoints using the `umount` command.

```
# umount /fsvol_may_23  
# umount /dev/vx/dsk/fsvol/vol1:may_23
```

Note: You do not need to run the `fsck` utility on Storage Checkpoint pseudo devices because pseudo devices are part of the actual file system.

Converting a data Storage Checkpoint to a nodata Storage Checkpoint

A nodata Storage Checkpoint does not contain actual file data. Instead, this type of Storage Checkpoint contains a collection of markers indicating the location of all the changed blocks since the Storage Checkpoint was created.

See [“Types of Storage Checkpoints”](#) on page 88.

You can use either the synchronous or asynchronous method to convert a data Storage Checkpoint to a nodata Storage Checkpoint; the asynchronous method is the default method. In a synchronous conversion, `fsckptadm` waits for all files to undergo the conversion process to “nodata” status before completing the operation. In an asynchronous conversion, `fsckptadm` returns immediately and marks the Storage Checkpoint as a nodata Storage Checkpoint even though the Storage Checkpoint's data blocks are not immediately returned to the pool of free blocks in the file system. The Storage Checkpoint deallocates all of its file data blocks in the background and eventually returns them to the pool of free blocks in the file system.

If all of the older Storage Checkpoints in a file system are nodata Storage Checkpoints, use the synchronous method to convert a data Storage Checkpoint to a nodata Storage Checkpoint. If an older data Storage Checkpoint exists in the

file system, use the asynchronous method to mark the Storage Checkpoint you want to convert for a delayed conversion. In this case, the actual conversion will continue to be delayed until the Storage Checkpoint becomes the oldest Storage Checkpoint in the file system, or all of the older Storage Checkpoints have been converted to nodata Storage Checkpoints.

Note: You cannot convert a nodata Storage Checkpoint to a data Storage Checkpoint because a nodata Storage Checkpoint only keeps track of the location of block changes and does not save the content of file data blocks.

Showing the difference between a data and a nodata Storage Checkpoint

The following example shows the difference between data Storage Checkpoints and nodata Storage Checkpoints.

Note: A nodata Storage Checkpoint does not contain actual file data.

See [“Converting a data Storage Checkpoint to a nodata Storage Checkpoint”](#) on page 94.

To show the difference between Storage Checkpoints

- 1 Create a file system and mount it on `/mnt0`, as in the following example:

```
# mkfs -F vxfs /dev/vx/rdisk/dg1/test0

version 7 layout
134217728 sectors, 67108864 blocks of size 1024, log \

size 65536 blocks, largefiles supported
# mount -F /dev/vx/rdisk/dg1/test0 /mnt0
```

- 2 Create a small file with a known content, as in the following example:

```
# echo "hello, world" > /mnt0/file
```

- 3 Create a Storage Checkpoint and mount it on `/mnt0@5_30pm`, as in the following example:

```
# fsckptadm create ckpt@5_30pm /mnt0
# mkdir /mnt0@5_30pm
# mount -F vxfs -o ckpt=ckpt@5_30pm \
/dev/vx/dsk/dg1/test0:ckpt@5_30pm /mnt0@5_30pm
```

- 4 Examine the content of the original file and the Storage Checkpoint file:

```
# cat /mnt0/file
hello, world
# cat /mnt0@5_30pm/file
hello, world
```

- 5 Change the content of the original file:

```
# echo "goodbye" > /mnt0/file
```

- 6 Examine the content of the original file and the Storage Checkpoint file. The original file contains the latest data while the Storage Checkpoint file still contains the data at the time of the Storage Checkpoint creation:

```
# cat /mnt0/file
goodbye
# cat /mnt0@5_30pm/file
hello, world
```

- 7 Unmount the Storage Checkpoint, convert the Storage Checkpoint to a nodata Storage Checkpoint, and mount the Storage Checkpoint again:

```
# umount /mnt0@5_30pm
# fsckptadm -s set nodata ckpt@5_30pm /mnt0
# mount -F vxfs -o ckpt=ckpt@5_30pm \
/dev/vx/dsk/dg1/test0:ckpt@5_30pm /mnt0@5_30pm
```

- 8 Examine the content of both files. The original file must contain the latest data:

```
# cat /mnt0/file
goodbye
```

You can traverse and read the directories of the nodata Storage Checkpoint; however, the files contain no data, only markers to indicate which block of the file has been changed since the Storage Checkpoint was created:

```
# ls -l /mnt0@5_30pm/file
-rw-r--r--  1 root      other 13 Jul 13 17:13 \
# cat /mnt0@5_30pm/file
cat: read error: No such file or directory
```


Converting multiple Storage Checkpoints

You can convert Storage Checkpoints to nodata Storage Checkpoints, when dealing with older Storage Checkpoints on the same file system.

To convert multiple Storage Checkpoints

- 1 Create a file system and mount it on /mnt0:

```
# mkfs -F vxfs /dev/vx/rdisk/dg1/test0
version 7 layout
13417728 sectors, 67108864 blocks of size 1024, log \
size 65536 blocks largefiles supported
# mount -F vxfs /dev/vx/dsk/dg1/test0 /mnt0
```

- 2 Create four data Storage Checkpoints on this file system, note the order of creation, and list them:

```
# fsckptadm create oldest /mnt0
# fsckptadm create older /mnt0
# fsckptadm create old /mnt0
# fsckptadm create latest /mnt0
# fsckptadm list /mnt0

/mnt0
latest:
  ctime           = Mon 26 Jul 11:56:55 2004
  mtime           = Mon 26 Jul 11:56:55 2004
  flags           = largefiles
old:
  ctime           = Mon 26 Jul 11:56:51 2004
  mtime           = Mon 26 Jul 11:56:51 2004
  flags           = largefiles
older:
  ctime           = Mon 26 Jul 11:56:46 2004
  mtime           = Mon 26 Jul 11:56:46 2004
  flags           = largefiles
oldest:
  ctime           = Mon 26 Jul 11:56:41 2004
  mtime           = Mon 26 Jul 11:56:41 2004
  flags           = largefiles
```

- 3 Try to convert synchronously the `latest` Storage Checkpoint to a `nodata` Storage Checkpoint. The attempt will fail because the Storage Checkpoints older than the `latest` Storage Checkpoint are data Storage Checkpoints, namely the Storage Checkpoints `old`, `older`, and `oldest`:

```
# fsckptadm -s set nodata latest /mnt0
UX:vxfs fsckptadm: ERROR: V-3-24632: Storage Checkpoint
set failed on latest. File exists (17)
```

- 4 You can instead convert the `latest` Storage Checkpoint to a `nodata` Storage Checkpoint in a delayed or asynchronous manner.

```
# fsckptadm set nodata latest /mnt0
```

- 5 List the Storage Checkpoints, as in the following example. You will see that the `latest` Storage Checkpoint is marked for conversion in the future.

```
# fsckptadm list /mnt0
/mnt0
latest:
  ctime          = Mon 26 Jul 11:56:55 2004
  mtime          = Mon 26 Jul 11:56:55
  flags          = nodata, largefiles, delayed
old:
  ctime          = Mon 26 Jul 11:56:51 2004
  mtime          = Mon 26 Jul 11:56:51 2004
  flags          = largefiles
older:
  ctime          = Mon 26 Jul 11:56:46 2004
  mtime          = Mon 26 Jul 11:56:46 2004
  flags          = largefiles
oldest:
  ctime          = Mon 26 Jul 11:56:41 2004
  mtime          = Mon 26 Jul 11:56:41 2004
  flags          = largefiles
```

Creating a delayed `nodata` Storage Checkpoint

You can combine the three previous steps and create the `latest` Storage Checkpoint as a `nodata` Storage Checkpoint. The creation process will detect the presence of the older data Storage Checkpoints and create the `latest` Storage Checkpoint as a delayed `nodata` Storage Checkpoint.

To create a delayed nodata Storage Checkpoint

1 Remove the latest Storage Checkpoint.

```
# fsckptadm remove latest /mnt0
# fsckptadm list /mnt0
/mnt0
old:
  ctime          = Mon 26 Jul 11:56:51 2004
  mtime          = Mon 26 Jul 11:56:51 2004
  flags          = largefiles
older:
  ctime          = Mon 26 Jul 11:56:46 2004
  mtime          = Mon 26 Jul 11:56:46 2004
  flags          = largefiles
oldest:
  ctime          = Mon 26 Jul 11:56:41 2004
  mtime          = Mon 26 Jul 11:56:41 2004
  flags          = largefiles
```

2 Recreate the latest Storage Checkpoint as a nodata Storage Checkpoint.

```
# fsckptadm -n create latest /mnt0
# fsckptadm list /mnt0
/mnt0
latest:
  ctime          = Mon 26 Jul 12:06:42 2004
  mtime          = Mon 26 Jul 12:06:42 2004
  flags          = nodata, largefiles, delayed
old:
  ctime          = Mon 26 Jul 11:56:51 2004
  mtime          = Mon 26 Jul 11:56:51 2004
  flags          = largefiles
older:
  ctime          = Mon 26 Jul 11:56:46 2004
  mtime          = Mon 26 Jul 11:56:46 2004
  flags          = largefiles
oldest:
  ctime          = Mon 26 Jul 11:56:41 2004
  mtime          = Mon 26 Jul 11:56:41 2004
  flags          = largefiles
```

- 3 Convert the `oldest` Storage Checkpoint to a `nodata` Storage Checkpoint because no older Storage Checkpoints exist that contain data in the file system.

Note: This step can be done synchronously.

```
# fsckptadm -s set nodata oldest /mnt0
# fsckptadm list /mnt0
/mnt0
latest:
  ctime          = Mon 26 Jul 12:06:42 2004
  mtime          = Mon 26 Jul 12:06:42 2004
  flags          = nodata, largefiles, delayed
old:
  ctime          = Mon 26 Jul 11:56:51 2004
  mtime          = Mon 26 Jul 11:56:51 2004
  flags          = largefiles
older:
  ctime          = Mon 26 Jul 11:56:46 2004
  mtime          = Mon 26 Jul 11:56:46 2004
  flags          = largefiles
oldest:
  ctime          = Mon 26 Jul 11:56:41 2004
  mtime          = Mon 26 Jul 11:56:41 2004
  flags          = nodata, largefiles
```

4 Remove the `older` and `old` Storage Checkpoints.

```
# fsckptadm remove older /mnt0
# fsckptadm remove old /mnt0
# fsckptadm list /mnt0
/mnt0
latest:
    ctime          = Mon 26 Jul 12:06:42 2004
    mtime          = Mon 26 Jul 12:06:42 2004
    flags          = nodata, largefiles
oldest:
    ctime          = Mon 26 Jul 11:56:41 2004
    mtime          = Mon 26 Jul 11:56:41 2004
    flags          = nodata, largefiles
```

Note: After you remove the `older` and `old` Storage Checkpoints, the `latest` Storage Checkpoint is automatically converted to a `nodata` Storage Checkpoint because the only remaining `older` Storage Checkpoint (`oldest`) is already a `nodata` Storage Checkpoint:

Space management considerations

Several operations, such as removing or overwriting a file, can fail when a file system containing Storage Checkpoints runs out of space. If the system cannot allocate sufficient space, the operation will fail.

Database applications usually preallocate storage for their files and may not expect a write operation to fail. If a file system runs out of space, the kernel automatically removes Storage Checkpoints and attempts to complete the write operation after sufficient space becomes available. The kernel removes Storage Checkpoints to prevent commands, such as `rm`, from failing under an out-of-space (`ENOSPC`) condition.

See the `rm(1M)` manual page.

When the kernel automatically removes the Storage Checkpoints, it applies the following policies:

- Remove as few Storage Checkpoints as possible to complete the operation.
- Never select a non-removable Storage Checkpoint.
- Select a `nodata` Storage Checkpoint only when data Storage Checkpoints no longer exist.

- Remove the oldest Storage Checkpoint first.

Restoring a file system from a Storage Checkpoint

Mountable data Storage Checkpoints on a consistent and undamaged file system can be used by backup and restore applications to restore either individual files or an entire file system. Restoration from Storage Checkpoints can also help recover incorrectly modified files, but typically cannot recover from hardware damage or other file system integrity problems.

Note: For hardware or other integrity problems, Storage Checkpoints must be supplemented by backups from other media.

Files can be restored by copying the entire file from a mounted Storage Checkpoint back to the primary fileset. To restore an entire file system, you can designate a mountable data Storage Checkpoint as the primary fileset using the `fsckpt_restore` command.

See the `fsckpt_restore(1M)` manual page.

When using the `fsckpt_restore` command to restore a file system from a Storage Checkpoint, all changes made to that file system after that Storage Checkpoint's creation date are permanently lost. The only Storage Checkpoints and data preserved are those that were created at the same time, or before, the selected Storage Checkpoint's creation. The file system cannot be mounted when `fsckpt_restore` is invoked.

Note: Files can be restored very efficiently by applications using the `fsckpt_fbmap(3)` library function to restore only modified portions of a file's data.

Restoring a file from a Storage Checkpoint

The following example restores a file, `MyFile.txt`, which resides in your home directory, from the Storage Checkpoint `CKPT1` to the device `/dev/vx/dsk/vol-01`. The mount point for the device is `/home`.

To restore a file from a Storage Checkpoint

- 1 Create the Storage Checkpoint CKPT1 of /home.

```
$ fckptadm create CKPT1 /home
```

- 2 Mount Storage Checkpoint CKPT1 on the directory /home/checkpoints/mar_4.

```
$ mount -F vxfs -o ckpt=CKPT1 /dev/vx/dsk/dg1/vol- \
01:CKPT1 /home/checkpoints/mar_4
```

- 3 Delete the file MyFile.txt from your home directory.

```
$ cd /home/users/me
$ rm MyFile.txt
```

- 4 Go to the /home/checkpoints/mar_4/users/me directory, which contains the image of your home directory.

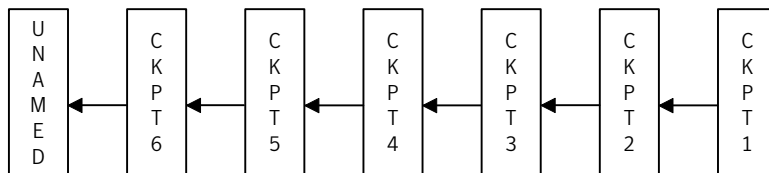
```
$ cd /home/checkpoints/mar_4/users/me
$ ls -l
-rw-r--r--  1 me  staff  14910   Mar 4   17:09  MyFile.txt
```

- 5 Copy the file MyFile.txt to your home directory.

```
$ cp MyFile.txt /home/users/me
$ cd /home/users/me
$ ls -l
-rw-r--r--  1 me  staff  14910   Mar 4   18:21  MyFile.txt
```

Restoring a file system from a Storage Checkpoint

The following example restores a file system from the Storage Checkpoint CKPT3. The filesets listed before the restoration show an unnamed root fileset and six Storage Checkpoints.



To restore a file system from a Storage Checkpoint

1 Run the `fsckpt_restore` command:

```
# fsckpt_restore -l /dev/vx/dsk/dg1/vol2
/dev/vx/dsk/dg1/vol2:
UNNAMED:
    ctime          = Thu 08 May 2004 06:28:26 PM PST
    mtime          = Thu 08 May 2004 06:28:26 PM PST
    flags          = largefiles, file system root
CKPT6:
    ctime          = Thu 08 May 2004 06:28:35 PM PST
    mtime          = Thu 08 May 2004 06:28:35 PM PST
    flags          = largefiles
CKPT5:
    ctime          = Thu 08 May 2004 06:28:34 PM PST
    mtime          = Thu 08 May 2004 06:28:34 PM PST
    flags          = largefiles, nomount
CKPT4:
    ctime          = Thu 08 May 2004 06:28:33 PM PST
    mtime          = Thu 08 May 2004 06:28:33 PM PST
    flags          = largefiles
CKPT3:
    ctime          = Thu 08 May 2004 06:28:36 PM PST
    mtime          = Thu 08 May 2004 06:28:36 PM PST
    flags          = largefiles
CKPT2:
    ctime          = Thu 08 May 2004 06:28:30 PM PST
    mtime          = Thu 08 May 2004 06:28:30 PM PST
    flags          = largefiles
CKPT1:
    ctime          = Thu 08 May 2004 06:28:29 PM PST
    mtime          = Thu 08 May 2004 06:28:29 PM PST
    flags          = nodata, largefiles
```

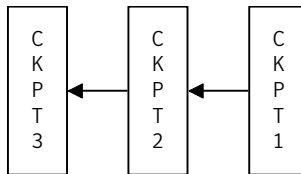

2 In this example, select the Storage Checkpoint CKPT3 as the new root filesset:

```
Select Storage Checkpoint for restore operation
or <Control/D> (EOF) to exit
or <Return> to list Storage Checkpoints: CKPT3
CKPT3:
    ctime          = Thu 08 May 2004 06:28:31 PM PST
    mtime          = Thu 08 May 2004 06:28:36 PM PST
    flags          = largefiles
UX:vxfs fsckpt_restore: WARNING: V-3-24640: Any file system
changes or Storage Checkpoints made after
Thu 08 May 2004 06:28:31 PM PST will be lost.
```

3 Type **y** to restore the file system from CKPT3:

```
Restore the file system from Storage Checkpoint CKPT3 ?
(ynq) y
(Yes)
UX:vxfs fscckpt_restore: INFO: V-3-23760: File system
restored from CKPT3
```

If the filesets are listed at this point, it shows that the former UNNAMED root fileset and CKPT6, CKPT5, and CKPT4 were removed, and that CKPT3 is now the primary fileset. CKPT3 is now the fileset that will be mounted by default.



4 Run the `fscckpt_restore` command:

```
# fscckpt_restore -l /dev/vx/dsk/dg1/vol2
/dev/vx/dsk/dg1/vol2:
CKPT3:
    ctime      = Thu 08 May 2004 06:28:31 PM PST
    mtime      = Thu 08 May 2004 06:28:36 PM PST
    flags      = largefiles, file system root
CKPT2:
    ctime      = Thu 08 May 2004 06:28:30 PM PST
    mtime      = Thu 08 May 2004 06:28:30 PM PST
    flags      = largefiles
CKPT1:
    ctime      = Thu 08 May 2004 06:28:29 PM PST
    mtime      = Thu 08 May 2004 06:28:29 PM PST
    flags      = nodata, largefiles
Select Storage Checkpoint for restore operation
or <Control/D> (EOF) to exit
or <Return> to list Storage Checkpoints:
```

Storage Checkpoint quotas

VxFS provides options to the `fscckptadm` command interface to administer Storage Checkpoint quotas. Storage Checkpoint quotas set the following limits on the amount of space used by all Storage Checkpoints of a primary file set:

hard limit	An absolute limit that cannot be exceeded. If a hard limit is exceeded, all further allocations on any of the Storage Checkpoints fail, but existing Storage Checkpoints are preserved.
soft limit	Must be lower than the hard limit. If a soft limit is exceeded, no new Storage Checkpoints can be created. The number of blocks used must return below the soft limit before more Storage Checkpoints can be created. An alert and console message are generated.

In case of a hard limit violation, various solutions are possible, enacted by specifying or not specifying the `-f` option for the `fscckptadm` utility.

See the `fscckptadm(1M)` manual page.

Specifying or not specifying the `-f` option has the following effects:

- If the `-f` option is not specified, one or many removable Storage Checkpoints are deleted to make space for the operation to succeed. This is the default solution.
- If the `-f` option is specified, all further allocations on any of the Storage Checkpoints fail, but existing Storage Checkpoints are preserved.

Note: Sometimes if a file is removed while it is opened by another process, the removal process is deferred until the last close. Because the removal of a file may trigger pushing data to a “downstream” Storage Checkpoint (that is, the next older Storage Checkpoint), a filesset hard limit quota violation may occur. In this scenario, the hard limit is relaxed to prevent an inode from being marked bad. This is also true for some asynchronous inode operations.

Quotas

This chapter includes the following topics:

- [About quota limits](#)
- [About quota files on Veritas File System](#)
- [About quota commands](#)
- [Using quotas](#)

About quota limits

Veritas File System (VxFS) supports user quotas. The quota system limits the use of two principal resources of a file system: files and data blocks. For each of these resources, you can assign quotas to individual users to limit their usage.

You can set the following kinds of limits for each of the two resources:

hard limit	An absolute limit that cannot be exceeded under any circumstances.
soft limit	Must be lower than the hard limit, and can be exceeded, but only for a limited time. The time limit can be configured on a per-file system basis only. The VxFS default limit is seven days.

Soft limits are typically used when a user must run an application that could generate large temporary files. In this case, you can allow the user to exceed the quota limit for a limited time. No allocations are allowed after the expiration of the time limit. Use the `edquota` command to set limits.

See [“Using quotas”](#) on page 110.

Although file and data block limits can be set individually for each user, the time limits apply to the file system as a whole. The quota limit information is associated with user IDs and is stored in a user quota file.

See [“About quota files on Veritas File System”](#) on page 110.

The quota soft limit can be exceeded when VxFS preallocates space to a file. Quota limits cannot exceed two terabytes on a Version 5 disk layout.

See [“About extent attributes”](#) on page 63.

About quota files on Veritas File System

A quotas file (named `quotas`) must exist in the root directory of a file system for any of the quota commands to work. The files in the file system's mount point are referred to as the external quotas file. VxFS also maintains an internal quotas file for its own use.

The quota administration commands read and write to the external quotas file to obtain or change usage limits. VxFS uses the internal file to maintain counts of data blocks and inodes used by each user. When quotas are turned on, the quota limits are copied from the external quotas file into the internal quotas file. While quotas are on, all the changes in the usage information and changes to quotas are registered in the internal quotas file. When quotas are turned off, the contents of the internal quotas file are copied into the external quotas file so that all data between the two files is synchronized.

About quota commands

Most of the quotas commands in VxFS are similar to BSD quotas commands. In general, quota administration for VxFS is performed using commands similar to HFS quota commands. The VxFS mount command supports a special mount option (`-o quotas`), that can be used to turn on quotas at mount time.

For additional information on the quota commands, see the corresponding manual pages. In general, quota administration for VxFS is similar to that of BSD quotas.

Note: When VxFS file systems are exported via NFS, the VxFS quota commands on the NFS client cannot query or edit quotas. You can use the VxFS quota commands on the server to query or edit quotas.

Using quotas

The VxFS quota commands are used to manipulate quotas.

Turning on quotas

To use the quota functionality on a file system, quotas must be turned on. You can turn quotas on at mount time or after a file system is mounted.

To turn on quotas

- ◆ To turn on user quotas for a VxFS file system, enter:

```
# quotaon /mount_point
```

Turning on quotas at mount time

Quotas can be turned on with the `mount` command when you mount a file system.

To turn on quotas at mount time

- ◆ To turn on user quotas for a file system at mount time, enter:

```
# mount -F vxfs -o quota special /mount_point
```

Editing user quotas

You can set up user quotas using the `edquota` command. You must have superuser privileges to edit quotas.

`edquota` creates a temporary file for the given user; this file contains on-disk quotas for each mounted file system that has a quotas file. It is not necessary that quotas be turned on for `edquota` to work. However, the quota limits are applicable only after quotas are turned on for a given file system.

To edit quotas

- ◆ Use the `edquota` command to edit quotas:

```
# edquota username
```

Modifying time limits

The soft and hard limits can be modified or assigned values. For any user, usage can never exceed the hard limit after quotas are turned on.

Modified time limits apply to the entire file system and cannot be set selectively for each user.

To modify time limits

- ◆ Specify the `-t` option to modify time limits for any user:

```
# edquota -t
```

Viewing disk quotas and usage

Use the `quota` command to view a user's disk quotas and usage on VxFS file systems.

To display disk quotas and usage

- ◆ To display a user's quotas and disk usage on all mounted VxFS file systems where the quotas file exists, enter:

```
# quota -v username
```

Displaying blocks owned by users or groups

Use the `quot` command to display the number of blocks owned by each user in a file system.

To display the number of blocks owned by users

- ◆ To display the number of files and the space owned by each user, enter:

```
# quot -f filesystem
```

Turning off quotas

Use the `vxquotaoff` command to turn off quotas.

To turn off quotas

- ◆ To turn off quotas for a mounted file system, enter:

```
# quotaoff /mount_point
```


File Change Log

This chapter includes the following topics:

- [About File Change Log](#)
- [About the File Change Log file](#)
- [File Change Log administrative interface](#)
- [File Change Log programmatic interface](#)
- [Summary of API functions](#)
- [Reverse path name lookup](#)

About File Change Log

The VxFS File Change Log (FCL) tracks changes to files and directories in a file system.

Applications that typically use the FCL are usually required to:

- scan an entire file system or a subset
- discover changes since the last scan

These applications may include: backup utilities, webcrawlers, search engines, and replication programs.

Note: The FCL tracks when the data has changed and records the change type, but does not track the actual data changes. It is the responsibility of the application to examine the files to determine the changed data.

FCL functionality is a separately licensable feature.

See the *Veritas Storage Foundation Release Notes*.

About the File Change Log file

File Change Log records file system changes such as creates, links, unlinks, renaming, data appended, data overwritten, data truncated, extended attribute modifications, holes punched, and miscellaneous file property updates.

Note: FCL is supported only on disk layout Version 6 and later.

FCL stores changes in a sparse file in the file system namespace. The FCL file is located in `mount_point/lost+found/changelog`. The FCL file behaves like a regular file, but some operations are prohibited. The standard system calls `open(2)`, `lseek(2)`, `read(2)` and `close(2)` can access the data in the FCL, while the `write(2)`, `mmap(2)` and `rename(2)` calls are not allowed.

Warning: Although some standard system calls are currently supported, the FCL file might be pulled out of the namespace in future VxFS release and these system calls may no longer work. It is recommended that all new applications be developed using the programmatic interface.

The FCL log file contains both the information about the FCL, which is stored in the FCL superblock, and the changes to files and directories in the file system, which is stored as FCL records.

See [“File Change Log programmatic interface”](#) on page 117.

In 4.1, the structure of the File Change Log file was exposed through the `/opt/VRTS/include/sys/fs/fcl.h` header file. In this release, the internal structure of the FCL file is opaque. The recommended mechanism to access the FCL is through the API described by the `/opt/VRTSfssdk/5.0/include/vxfsutil.h` header file.

The `/opt/VRTS/include/sys/fs/fcl.h` header file is included in this release to ensure that applications accessing the FCL with the 4.1 header file do not break. New applications should use the new FCL API described in `/opt/VRTSfssdk/5.0/include/vxfsutil.h`. Existing applications should also be modified to use the new FCL API.

With the addition of new record types, the FCL version in this release has been updated to 4. To provide backward compatibility for the existing applications, this release supports multiple FCL versions. Users now have the flexibility of specifying the FCL version for new FCLs. The default FCL version is 4.

See the `fcladm(1M)` man page.

File Change Log administrative interface

The FCL can be set up and tuned through the `fcladm` and `vxtunefs` VxFS administrative commands.

See the `fcladm(1M)` and `vxtunefs(1M)` manual pages.

The FCL keywords for `fcladm` are as follows:

<code>clear</code>	Disables the recording of the audit, open, close, and statistical events after it has been set.
<code>dump</code>	Creates a regular file image of the FCL file that can be downloaded too an off-host processing system. This file has a different format than the FCL file.
<code>on</code>	Activates the FCL on a mounted file system. VxFS 5.0 supports either FCL Versions 3 or 4. If no version is specified, the default is Version 4. Use <code>fcladm on</code> to specify the version.
<code>print</code>	Prints the contents of the FCL file starting from the specified offset.
<code>restore</code>	Restores the FCL file from the regular file image of the FCL file created by the <code>dump</code> keyword.
<code>rm</code>	Removes the FCL file. You must first deactivate the FCL with the <code>off</code> keyword, before you can remove the FCL file.
<code>set</code>	Enables the recording of events specified by the 'eventlist' option. See the <code>fcladm(1M)</code> manual page.
<code>state</code>	Writes the current state of the FCL to the standard output.
<code>sync</code>	Brings the FCL to a stable state by flushing the associated data of an FCL recording interval.

The FCL tunable parameters for `vxtunefs` are as follows:

<code>fcl_keeptime</code>	Specifies the duration in seconds that FCL records stay in the FCL file before they can be purged. The first records to be purged are the oldest ones, which are located at the beginning of the file. Additionally, records at the beginning of the file can be purged if allocation to the FCL file exceeds <code>fcl_maxalloc</code> bytes. The default value of <code>fcl_keeptime</code> is 0. If the <code>fcl_maxalloc</code> parameter is set, records are purged from the FCL file if the amount of space allocated to the FCL file exceeds <code>fcl_maxalloc</code> . This is true even if the elapsed time the records have been in the log is less than the value of <code>fcl_keeptime</code> .
<code>fcl_maxalloc</code>	Specifies the maximum number of spaces in bytes to be allocated to the FCL file. When the space allocated exceeds <code>fcl_maxalloc</code> , a hole is punched at the beginning of the file. As a result, records are purged and the first valid offset (<code>fc_off</code>) is updated. In addition, <code>fcl_maxalloc</code> may be violated if the oldest record has not reached <code>fcl_keeptime</code> . The minimum value of <code>fcl_maxalloc</code> is 4 MB. The default value is <code>fs_size/33</code> .
<code>fcl_winterval</code>	Specifies the time in seconds that must elapse before the FCL records an overwrite, extending write, or a truncate. This helps to reduce the number of repetitive records in the FCL. The <code>fcl_winterval</code> timeout is per inode. If an inode happens to go out of cache and returns, its write interval is reset. As a result, there could be more than one write record for that file in the same write interval. The default value is 3600 seconds.
<code>fcl_ointerval</code>	The time interval in seconds within which subsequent opens of a file do not produce an additional FCL record. This helps to reduce the number of repetitive records logged in the FCL file. If the tracking of access information is also enabled, a subsequent file open even within the <code>fcl_ointerval</code> may produce a record, if it is opened by a different user. Similarly, if the inode is bumped out of cache, this may also produce more than one record within the same open interval. The default value is 600 sec.

Either or both `fcl_maxalloc` and `fcl_keeptime` must be set to activate the FCL feature. The following are examples of using the `fcladm` command.

To activate FCL for a mounted file system, type the following:

```
# fcladm on mount_point
```

To deactivate the FCL for a mounted file system, type the following:

```
# fcladm off mount_point
```

To remove the FCL file for a mounted file system, on which FCL must be turned off, type the following:

```
# fcladm rm mount_point
```

To obtain the current FCL state for a mounted file system, type the following:

```
# fcladm state mount_point
```

To enable tracking of the file opens along with access information with each event in the FCL, type the following:

```
# fcladm set fileopen,accessinfo mount_point
```

To stop tracking file I/O statistics in the FCL, type the following:

```
# fcladm clear filestats mount_point
```

Print the on-disk FCL super-block in text format to obtain information about the FCL file by using offset 0. Because the FCL on-disk super-block occupies the first block of the FCL file, the first and last valid offsets into the FCL file can be determined by reading the FCL super-block and checking the `fc_off` field. Enter:

```
# fcladm print 0 mount_point
```

To print the contents of the FCL in text format, of which the offset used must be 32-byte aligned, enter:

```
# fcladm print offset mount_point
```

File Change Log programmatic interface

VxFS provides an enhanced API to simplify reading and parsing the FCL file in two ways:

Simplified reading	The API simplifies user tasks by reducing additional code needed to parse FCL file entries. In 4.1, to obtain event information such as a remove or link, the user was required to write additional code to get the name of the removed or linked file. In this release, the API allows the user to directly read an assembled record. The API also allows the user to specify a filter to indicate a subset of the event records of interest.
--------------------	--

Backward compatibility	Providing API access for the FCL feature allows backward compatibility for applications. The API allows applications to parse the FCL file independent of the FCL layout changes. Even if the hidden disk layout of the FCL changes, the API automatically translates the returned data to match the expected output record. As a result, the user does not need to modify or recompile the application due to changes in the on-disk FCL layout.
------------------------	---

See [“Reverse path name lookup”](#) on page 120.

The following sample code fragment reads the FCL superblock, checks that the state of the FCL is `VX_FCLS_ON`, issues a call to `vxfs_fcl_sync` to obtain a finishing offset to read to, determines the first valid offset in the FCL file, then reads the entries in 8K chunks from this offset. The section process fcl entries is what an application developer must supply to process the entries in the FCL file.

```
#include <stdint.h>
#include <stdio.h>
#include <stdlib.h>
#include <sys/types.h>
#include <sys/fcntl.h>
#include <errno.h>
#include <fcl.h>
#include <vxfsutil.h>
#define FCL_READSZ 8192
char* fclname = "/mnt/lost+found/changelog";
int read_fcl(char* fclname)
{
    struct fcl_sb fclsb;
    uint64_t off, lastoff;
    size_t size;
    char buf[FCL_READSZ], *bufp = buf;
    int fd;
    int err = 0;
    if ((fd = open(fclname, O_RDONLY)) < 0) {
        return ENOENT;
    }
    if ((off = lseek(fd, 0, SEEK_SET)) != 0) {
        close(fd);
        return EIO;
    }
    size = read(fd, &fclsb, sizeof (struct fcl_sb));
    if (size < 0) {
        close(fd);
    }
}
```

```

        return EIO;
    }
    if (fclsb.fc_state == VX_FCLS_OFF) {
        close(fd);
        return 0;
    }
    if (err = vxfs_fcl_sync(fclname, &lastoff)) {
        close(fd);
        return err;
    }
    if ((off = lseek(fd, off_t, uint64_t)) != uint64_t) {
        close(fd);
        return EIO;
    }
    while (off < lastoff) {
        if ((size = read(fd, bufp, FCL_READSZ)) <= 0) {
            close(fd);
            return errno;
        }
        /* process fcl entries */
        off += size;
    }
    close(fd);
    return 0;
}

```

Summary of API functions

The following is a brief summary of File Change Log API functions:

<code>vxfs_fcl_close()</code>	Closes the FCL file and cleans up resources associated with the handle.
<code>vxfs_fcl_cookie()</code>	Returns an opaque structure that embeds the current FCL activation time and the current offset. This cookie can be saved and later passed to <code>vxfs_fcl_seek()</code> function to continue reading from where the application last stopped.
<code>vxfs_fcl_getinfo()</code>	Returns information such as the state and version of the FCL file.
<code>vxfs_fcl_open()</code>	Opens the FCL file and returns a handle that can be used for further operations.
<code>vxfs_fcl_read()</code>	Reads FCL records of interest into a buffer specified by the user.

`vxfs_fcl_seek()` Extracts data from the specified cookie and then seeks to the specified offset.

`vxfs_fcl_seektime()` Seeks to the first record in the FCL after the specified time.

Reverse path name lookup

The reverse path name lookup feature obtains the full path name of a file or directory from the inode number of that file or directory. The inode number is provided as an argument to the `vxlsino` administrative command, or the `vxfs_inotopath_gen(3)` application programming interface library function.

The reverse path name lookup feature can be useful for a variety of applications, such as for clients of the VxFS File Change Log feature, in backup and restore utilities, and for replication products. Typically, these applications store information by inode numbers because a path name for a file or directory can be very long, thus the need for an easy method of obtaining a path name.

An inode is a unique identification number for each file in a file system. An inode contains the data and metadata associated with that file, but does not include the file name to which the inode corresponds. It is therefore relatively difficult to determine the name of a file from an inode number. The `ncheck` command provides a mechanism for obtaining a file name from an inode identifier by scanning each directory in the file system, but this process can take a long period of time. The VxFS reverse path name lookup feature obtains path names relatively quickly.

Note: Because symbolic links do not constitute a path to the file, the reverse path name lookup feature cannot track symbolic links to files.

Because of the possibility of errors with processes renaming or unlinking and creating new files, it is advisable to perform a `lookup` or `open` with the path name and verify that the inode number matches the path names obtained.

See the `vxlsino(1M)`, `vxfs_inotopath_gen(3)`, and `vxfs_inotopath(3)` manual pages.

Multi-volume file systems

This chapter includes the following topics:

- [About multi-volume support](#)
- [About volume types](#)
- [Features implemented using multi-volume support](#)
- [About volume sets](#)
- [Creating multi-volume file systems](#)
- [Converting a single volume file system to a multi-volume file system](#)
- [Removing a volume from a multi-volume file system](#)
- [About allocation policies](#)
- [Assigning allocation policies](#)
- [Querying allocation policies](#)
- [Assigning pattern tables to directories](#)
- [Assigning pattern tables to file systems](#)
- [Allocating data](#)
- [Volume encapsulation](#)
- [Reporting file extents](#)
- [Load balancing](#)
- [Converting a multi-volume file system to a single volume file system](#)

About multi-volume support

VxFS provides support for multi-volume file systems when used in conjunction with the Veritas Volume Manager. Using multi-volume support (MVS), a single file system can be created over multiple volumes, each volume having its own properties. For example, it is possible to place metadata on mirrored storage while placing file data on better-performing volume types such as RAID-1+0 (striped and mirrored).

The MVS feature also allows file systems to reside on different classes of devices, so that a file system can be supported from both inexpensive disks and from expensive arrays. Using the MVS administrative interface, you can control which data goes on which volume types.

Note: MVS is available only on file systems using disk layout Version 6 or later. See [“About disk layouts”](#) on page 255.

About volume types

VxFS utilizes two types of volumes, one of which contains only data, referred to as `dataonly`, and the other of which can contain metadata or data, referred to as `metadataok`.

Data refers to direct extents, which contain user data, of regular files and named data streams in a file system.

Metadata refers to all extents that are not regular file or name data stream extents. This includes certain files that appear to be regular files, but are not, such as the File Change Log file.

A volume availability flag is set to specify if a volume is `dataonly` or `metadataok`. The volume availability flag can be set, cleared, and listed with the `fsvoladm` command.

See the `fsvoladm(1M)` manual page.

Features implemented using multi-volume support

The following features can be implemented using multi-volume support:

- Controlling where files are stored can be selected at multiple levels so that specific files or file hierarchies can be assigned to different volumes. This functionality is available in the Veritas File System Dynamic Storage Tiering (DST) feature.

See [“About Dynamic Storage Tiering”](#) on page 143.

- Placing the VxFS intent log on its own volume to minimize disk head movement and thereby increase performance. This functionality can be used to migrate from the Veritas QuickLog™ feature.
- Separating Storage Checkpoints so that data allocated to a Storage Checkpoint is isolated from the rest of the file system.
- Separating metadata from file data.
- Encapsulating volumes so that a volume appears in the file system as a file. This is particularly useful for databases that are running on raw volumes.
- Guaranteeing that a dataonly volume being unavailable does not cause a metadataok volume to be unavailable.

To use the multi-volume file system features, Veritas Volume Manager must be installed and the volume set feature must be accessible.

Volume availability

MVS guarantees that a dataonly volume being unavailable does not cause a metadataok volume to be unavailable. This allows you to mount a multi-volume file system even if one or more component dataonly volumes are missing.

The volumes are separated by whether metadata is allowed on the volume. An I/O error on a dataonly volume does not affect access to any other volumes. All VxFS operations that do not access the missing dataonly volume function normally.

Some VxFS operations that do not access the missing dataonly volume and function normally include the following:

- Mounting the multi-volume file system, regardless if the file system is read-only or read/write.
- Kernel operations.
- Performing a `fsck` replay. Logged writes are converted to normal writes if the corresponding volume is dataonly.
- Performing a full `fsck`.
- Using all other commands that do not access data on a missing volume.

Some operations that could fail if a dataonly volume is missing include:

- Reading or writing file data if the file's data extents were allocated from the missing dataonly volume.
- Using the `vxdump` command.

Volume availability is supported only on a file system with disk layout Version 7 or later.

Note: Do not mount a multi-volume system with the `ioerror=disable` or `ioerror=wdisable` mount options if the volumes have different availability properties. Symantec recommends the `ioerror=mdisable` mount option both for cluster mounts and for local mounts.

About volume sets

Veritas Volume Manager exports a data object called a volume set. A volume set is a container for one or more volumes, each of which can have its own geometry. Unlike the traditional Volume Manager volume, which can be used for raw I/O access or to contain a file system, a volume set can only be used to contain a VxFS file system.

The Volume Manager `vxvset` command is used to create and manage volume sets. Volume sets cannot be empty. When the last entry is removed, the volume set itself is removed.

Creating and managing volume sets

The following command examples show how to create and manage volume sets.

To create and manage volume sets

- 1 Create a new volume set from `vol1`:

```
# vxassist -g dg1 make vol1 10m
# vxvset -g dg1 make myvset vol1
```

- 2 Create two new volumes and add them to the volume set:

```
# vxassist -g dg1 make vol2 50m
# vxassist -g dg1 make vol3 50m
# vxvset -g dg1 addvol myvset vol2
# vxvset -g dg1 addvol myvset vol3
```

- 3 List the component volumes of the previously created volume set:

```
# vxvset -g dg1 list myvset
```

VOLUME	INDEX	LENGTH	STATE	CONTEXT
vol1	0	20480	ACTIVE	-
vol2	1	102400	ACTIVE	-
vol3	2	102400	ACTIVE	-

- 4 Use the `ls` command to see that when a volume set is created, the volumes contained by the volume set are removed from the namespace and are instead accessed through the volume set name:

```
# ls -l /dev/vx/rdisk/rootdg/myvset
```

crw-----	1	root	root	108,70009	May 21 15:37	/dev/vx/rdisk/rootdg/myvset
----------	---	------	------	-----------	--------------	-----------------------------

- 5 Create a volume, add it to the volume set, and use the `ls` command to see that when a volume is added to the volume set, it is no longer visible in the namespace:

```
# vxassist -g dg1 make vol4 50m
# ls -l /dev/vx/rdisk/rootdg/vol4
```

crw-----	1	root	root	108,70012	May 21 15:43	/dev/vx/rdisk/rootdg/vol4
----------	---	------	------	-----------	--------------	---------------------------

```
# vxvset -g dg1 addvol myvset vol4
# ls -l /dev/vx/rdisk/rootdg/vol4
```

/dev/vx/rdisk/rootdg/vol4: No such file or directory

Creating multi-volume file systems

When a multi-volume file system is created, all volumes are `dataonly`, except volume zero. The volume availability flag of volume zero cannot be set to `dataonly`.

As metadata cannot be allocated from `dataonly` volumes, enough metadata space should be allocated using `metadataok` volumes. The "file system out of space" error occurs if there is insufficient metadata space available, even if the `df` command shows that there is free space in the file system. The `fsvoladm` command can be used to see the free space in each volume and set the availability flag of the volume.

Unless otherwise specified, VxFS commands function the same on multi-volume file systems the same as the commands do on single-volume file systems.

Example of creating a multi-volume file system

The following procedure is an example of creating a multi-volume file system.

To create a multi-volume file system

- 1 After a volume set is created, create a VxFS file system by specifying the volume set name as an argument to `mkfs`:

```
# mkfs -F vxfs /dev/vx/rdisk/rootdg/myvset
version 7 layout
327680 sectors, 163840 blocks of size 1024,
log size 1024 blocks largefiles supported
```

After the file system is created, VxFS allocates space from the different volumes within the volume set.

- 2 List the component volumes of the volume set using of the `fsvoladm` command:

```
# mount -F vxfs /dev/vx/dsk/rootdg/myvset /mnt1
# fsvoladm list /mnt1
```

devid	size	used	avail	name
0	10240	1280	8960	vol1
1	51200	16	51184	vol2
2	51200	16	51184	vol3
3	51200	16	51184	vol4

- 3 Add a new volume by adding the volume to the volume set, then adding the volume to the file system:

```
# vxassist -g dgl make vol5 50m
# vxvset -g dgl addvol myvset vol5
# fsvoladm add /mnt1 vol5 50m
# fsvoladm list /mnt1
```

devid	size	used	avail	name
0	10240	1300	8940	vol1
1	51200	16	51184	vol2
2	51200	16	51184	vol3
3	51200	16	51184	vol4
4	51200	16	51184	vol5

- 4 List the volume availability flags using the `fsvoladm` command:

```
# fsvoladm queryflags /mnt1
volname      flags
vol1         metadataok
vol2         dataonly
vol3         dataonly
vol4         dataonly
vol5         dataonly
```

- 5 Increase the metadata space in the file system using the `fsvoladm` command:

```
# fsvoladm clearflags dataonly /mnt1 vol2
# fsvoladm queryflags /mnt1
volname      flags
vol1         metadataok
vol2         metadataok
vol3         dataonly
vol4         dataonly
vol5         dataonly
```

Converting a single volume file system to a multi-volume file system

The following procedure converts a traditional, single volume file system, `/mnt1`, on a single volume `vol1` in the diskgroup `dg1` to a multi-volume file system.

To convert a single volume file system

- 1 Determine the version of the volume's diskgroup:

```
# vxdg list dg1 | grep version: | awk '{ print $2 }'
105
```

- 2 If the version is less than 110, upgrade the diskgroup:

```
# vxdg upgrade dg1
```

- 3 Determine the disk layout version of the file system:

```
# vxupgrade /mnt1
Version 4
```

- 4 If the disk layout version is less than 6, upgrade to Version 7. For example, if the disk layout version is 4, upgrade to Version 7 as follows:

```
# vxupgrade -n 5 /mnt1
# vxupgrade -n 6 /mnt1
# vxupgrade -n 7 /mnt1
```

- 5 Unmount the file system:

```
# umount /mnt1
```

- 6 Convert the volume into a volume set:

```
# vxvset -g dg1 make vset1 vol1
```

- 7 Edit the `/etc/fstab` file to replace the volume device name, `vol1`, with the volume set name, `vset1`.

- 8 Mount the file system:

```
# mount -F vxfs /dev/vx/dsk/dg1/vset1 /mnt1
```

- 9 As necessary, create and add volumes to the volume set:

```
# vxassist -g dg1 make vol2 256M
# vxvset -g dg1 addvol vset1 vol2
```

- 10 Set the placement class tags on all volumes that do not have a tag:

```
# vxassist -g dg1 settag vol1 vxfs.placement_class.tier1
# vxassist -g dg1 settag vol2 vxfs.placement_class.tier2
```

- 11 Add the new volumes to the file system:

```
# fsvoladm add /mnt1 vol2 256m
```

Removing a volume from a multi-volume file system

Use the `fsvoladm remove` command to remove a volume from a multi-volume file system. The `fsvoladm remove` command fails if the volume being removed is the only volume in any allocation policy.

Forcibly removing a volume

If you must forcibly remove a volume from a file system, such as if a volume is permanently destroyed and you want to clean up the dangling pointers to the lost volume, use the `fsck -o zapvol=volname` command. The `zapvol` option performs a full file system check and zaps all inodes that refer to the specified volume. The `fsck` command prints the inode numbers of all files that the command destroys; the file names are not printed. The `zapvol` option only affects regular files if used on a `dataonly` volume. However, it could destroy structural files if used on a `metadataok` volume, which can make the file system unrecoverable. Therefore, the `zapvol` option should be used with caution on `metadataok` volumes.

Moving volume 0

Volume 0 in a multi-volume file system cannot be removed from the file system, but you can move volume 0 to different storage using the `vxassist move` command. The `vxassist` command creates any necessary temporary mirrors and cleans up the mirrors at the end of the operation.

To move volume 0

- ◆ Move volume 0:

```
# vxassist -g mydg move vol1 !mydg
```

About allocation policies

To make full use of the multi-volume support feature, you can create allocation policies that allow files or groups of files to be assigned to specified volumes within the volume set.

A policy specifies a list of volumes and the order in which to attempt allocations. A policy can be assigned to a file, a file system, or a Storage Checkpoint created from a file system. When policies are assigned to objects in the file system, you must specify how the policy maps to both metadata and file data. For example, if a policy is assigned to a single file, the file system must know where to place both the file data and metadata. If no policies are specified, the file system places data randomly.

Assigning allocation policies

The following example shows how to assign allocation policies. The example volume set contains four volumes from different classes of storage.

To assign allocation policies

1 List the volumes in the volume set:

```
# vxvset -g rootdg list myvset
```

VOLUME	INDEX	LENGTH	STATE	CONTEXT
vol1	0	102400	ACTIVE	-
vol2	1	102400	ACTIVE	-
vol3	2	102400	ACTIVE	-
vol4	3	102400	ACTIVE	-

2 Create a file system on the myvset volume set and mount the file system:

```
# mkfs -F vxfs /dev/vx/rdisk/rootdg/myvset
version 7 layout
204800 sectors, 102400 blocks of size 1024,
log size 1024 blocks
largefiles supported
# mount -F vxfs /dev/vx/dsk/rootdg/myvset /mnt1
```

- 3 Define three allocation policies, `v1`, `bal_34`, and `rr_all`, that allocate from the volumes using different methods:

```
# fsapadm define /mnt1 v1 vol1
# fsapadm define -o balance -c 64k /mnt1 bal_34 vol3 vol4
# fsapadm define -o round-robin /mnt1 rr_all vol1 vol2 vol3 vol4
# fsapadm list /mnt1
```

name	order	flags	chunk	num	comps
rr_all	round-robin	0	0	4	vol1, vol2, vol3, vol4
bal_34	balance	0	64.000KB	2	vol3, vol4
v1	as-given	0	0	1	vol1

These policies allocate from the volumes as follows:

<code>v1</code>	Allocations are restricted to <code>vol1</code> .
<code>bal_34</code>	Balanced allocations between <code>vol3</code> and <code>vol4</code> .
<code>rr_all</code>	Round-robin allocations from all four volumes.

- 4 Assign the policies to various objects in the file system. The data policy must be specified before the metadata policy:

```
# fsapadm assignfile -f inherit /mnt1/appdir bal_34 v1
# fsapadm assignfs /mnt1 rr_all ''
```

These assignments will cause allocations for any files below `/mnt1/appdir` to use `bal_34` for data, and `v1` for metadata. Allocations for other files in the file system will default to the file system-wide policies given in `assignfs`, with data allocations from `rr_all`, and metadata allocations using the default policy as indicated by the empty string (`''`). The default policy is as-given allocations from all metadata-eligible volumes.

Querying allocation policies

Querying an allocation policy displays the definition of the allocation policy.

The following example shows how to query allocation policies. The example volume set contains two volumes from different classes of storage.

To query allocation policies

- ◆ Query the allocation policy:

```
# fsapadm query /mnt1 bal_34
```

Assigning pattern tables to directories

A pattern table contains patterns against which a file's name and creating process' UID and GID are matched as a file is created in a specified directory. The first successful match is used to set the allocation policies of the file, taking precedence over inheriting per-file allocation policies.

See the `fsapadm(1M)` manual page.

The following example shows how to assign pattern tables to a directory in a volume set that contains two volumes from different classes of storage. The pattern table matches all files created in the directory `dir1` with the `.mp3` extension for any user or group ID and assigns the `mp3data` data policy and `mp3meta` metadata policy.

To assign pattern tables to directories

- 1 Define two allocation policies called `mp3data` and `mp3meta` to refer to the `vol1` and `vol2` volumes:

```
# fsapadm define /mnt1 mp3data vol1
# fsapadm define /mnt1 mp3meta vol2
```

- 2 Assign the pattern table:

```
# fsapadm assignfilepat dir1 *.mp3//mp3data/mp3meta/
```

Assigning pattern tables to file systems

A pattern table contains patterns against which a file's name and creating process' UID and GID are matched as a file is created in a directory. If the directory does not have its own pattern table or an inheritable allocation policy, the file system's pattern table takes effect. The first successful match is used to set the allocation policies of the file.

See the `fsapadm(1M)` manual page.

The following example shows how to assign pattern tables to a file system in a volume set that contains two volumes from different classes of storage. The pattern table is contained within the pattern file `mypatternfile`.

To assign pattern tables to directories

- 1 Define two allocation policies called `mydata` and `mymeta` to refer to the `vol1` and `vol2` volumes:

```
# fsapadm define /mnt1 mydata vol1
# fsapadm define /mnt1 mymeta vol2
```

- 2 Assign the pattern table:

```
# fsapadm assignfspat -F mypatternfile /mnt1
```

Allocating data

The following script creates a large number of files to demonstrate the benefit of allocating data:

```
i=1
while [ $i -lt 1000 ]
do
    dd if=/dev/zero of=/mnt1/$i bs=65536 count=1
    i=`expr $i + 1`
done
```

Before the script completes, `vol1` runs out of space even though space is still available on the `vol2` volume:

```
# fsvoladm list /mnt1
```

devid	size	used	avail	name
0	51200	51200	0	vol1
1	51200	221	50979	vol2

One possible solution is to define and assign an allocation policy that allocates user data to the least full volume.

You must have system administrator privileges to create, remove, or change policies, or to set file system or Storage Checkpoint level policies. Users can assign a pre-existing policy to their files if the policy allows that.

Policies can be inherited for new files. A file will inherit the allocation policy of the directory in which it resides if you run the `fsapadm assignfile -f inherit` command on the directory.

The following example defines an allocation policy that allocates data to the least full volume.

Allocating data from vol1 to vol2

- 1 Define an allocation policy, `lf_12`, that allocates user data to the least full volume between `vol1` and `vol2`:

```
# fsapadm define -o least-full /mnt1 lf_12 vol1 vol2
```

- 2 Assign the allocation policy `lf_12` as the data allocation policy to the file system mounted at `/mnt1`:

```
# fsapadm assignfs /mnt1 lf_12 ''
```

Metadata allocations use the default policy, as indicated by the empty string (`''`). The default policy is as-given allocations from all metadata-eligible volumes.

Volume encapsulation

Multi-volume support enables the ability to encapsulate an existing raw volume and make the volume contents appear as a file in the file system.

Encapsulating a volume involves the following actions:

- Adding the volume to an existing volume set.
- Adding the volume to the file system using `fsvoladm`.

Encapsulating a volume

The following example illustrates how to encapsulate a volume.

To encapsulate a volume

1 List the volumes:

```
# vxvset -g dg1 list myvset
VOLUME    INDEX    LENGTH    STATE    CONTEXT
vol1       0         102400    ACTIVE   -
vol2       1         102400    ACTIVE   -
```

The volume set has two volumes.

2 Create a third volume and copy the passwd file to the third volume:

```
# vxassist -g dg1 make dbvol 100m
# dd if=/etc/passwd of=/dev/vx/rdisk/rootdg/dbvol count=1
1+0 records in
1+0 records out
```

The third volume will be used to demonstrate how the volume can be accessed as a file, as shown later.

3 Create a file system on the volume set:

```
# mkfs -F vxfs /dev/vx/rdisk/rootdg/myvset
version 7 layout
204800 sectors, 102400 blocks of size 1024,
log size 1024 blocks
largefiles supported
```

4 Mount the volume set:

```
# mount -F vxfs /dev/vx/dsk/rootdg/myvset /mnt1
```

5 Add the new volume to the volume set:

```
# vxvset -g dg1 addvol myvset dbvol
```

6 Encapsulate dbvol:

```
# fsvoladm encapsulate /mnt1/dbfile dbvol 100m
# ls -l /mnt1/dbfile
-rw----- 1 root other 104857600 May 22 11:30 /mnt1/dbfile
```

7 Examine the contents of dbfile to see that it can be accessed as a file:

```
# head -2 /mnt1/dbfile
root:x:0:1:Super-User:/:/sbin/sh
daemon:x:1:1:/::
```

The passwd file that was written to the raw volume is now visible in the new file.

Note: If the encapsulated file is changed in any way, such as if the file is extended, truncated, or moved with an allocation policy or resized volume, or the volume is encapsulated with a bias, the file cannot be de-encapsulated.

Deencapsulating a volume

The following example illustrates how to deencapsulate a volume.

To deencapsulate a volume**1 List the volumes:**

```
# vxvset -g dgl list myvset
VOLUME   INDEX   LENGTH   STATE   CONTEXT
vol1      0      102400   ACTIVE  -
vol2      1      102400   ACTIVE  -
dbvol     2      102400   ACTIVE  -
```

The volume set has three volumes.

2 Deencapsulate dbvol:

```
# fsvoladm deencapsulate /mnt1/dbfile
```

Reporting file extents

MVS feature provides the capability for file-to-volume mapping and volume-to-file mapping via the `fsmap` and `fsvmap` commands. The `fsmap` command reports the

volume name, logical offset, and size of data extents, or the volume name and size of indirect extents associated with a file on a multi-volume file system. The `fsvmap` command maps volumes to the files that have extents on those volumes.

See the `fsmmap(1M)` and `fsvmap(1M)` manual pages.

The `fsmmap` command requires `open()` permission for each file or directory specified. Root permission is required to report the list of files with extents on a particular volume.

Examples of reporting file extents

The following examples show typical uses of the `fsmmap` and `fsvmap` commands.

Using the `fsmmap` command

- ◆ Use the `find` command to descend directories recursively and run `fsmmap` on the list of files:

```
# find . | fsmmap -
Volume  Extent Type      File
vol2    Data                ./file1
vol1    Data                ./file2
```

Using the `fsvmap` command

- 1 Report the extents of files on multiple volumes:

```
# fsvmap /dev/vx/rdsk/fstest/testvset vol1 vol2
vol1  /.
vol1  /ns2
vol1  /ns3
vol1  /file1
vol2  /file1
vol2  /file2
```

- 2 Report the extents of files that have either data or metadata on a single volume in all Storage Checkpoints, and indicate if the volume has file system metadata:

```
# fsvmap -mvC /dev/vx/rdsk/fstest/testvset vol1
Meta  Structural  vol1  //volume has filesystem metadata//
Data  UNNAMED       vol1  /.
Data  UNNAMED       vol1  /ns2
Data  UNNAMED       vol1  /ns3
Data  UNNAMED       vol1  /file1
Meta  UNNAMED       vol1  /file1
```

Load balancing

An allocation policy with the balance allocation order can be defined and assigned to files that must have their allocations distributed at random between a set of specified volumes. Each extent associated with these files are limited to a maximum size that is defined as the required chunk size in the allocation policy. The distribution of the extents is mostly equal if none of the volumes are full or disabled.

Load balancing allocation policies can be assigned to individual files or for all files in the file system. Although intended for balancing data extents across volumes, a load balancing policy can be assigned as a metadata policy if desired, without any restrictions.

Note: If a file has both a fixed extent size set and an allocation policy for load balancing, certain behavior can be expected. If the chunk size in the allocation policy is greater than the fixed extent size, all extents for the file are limited by the chunk size. For example, if the chunk size is 16 MB and the fixed extent size is 3 MB, then the largest extent that satisfies both the conditions is 15 MB. If the fixed extent size is larger than the chunk size, all extents are limited to the fixed extent size. For example, if the chunk size is 2 MB and the fixed extent size is 3 MB, then all extents for the file are limited to 3 MB.

Defining and assigning a load balancing allocation policy

The following example defines a load balancing policy and assigns the policy to the file, `/mnt/file.db`.

To define and assign the policy

- 1 Define the policy by specifying the `-o balance` and `-c` options:

```
# fsapadm define -o balance -c 2m /mnt loadbal vol1 vol2 vol3 vol4
```

- 2 Assign the policy:

```
# fsapadm assign /mnt/filedb loadbal meta
```

Rebalancing extents

Extents can be rebalanced by strictly enforcing the allocation policy. Rebalancing is generally required when volumes are added or removed from the policy or when the chunk size is modified. When volumes are removed from the volume set, any

extents on the volumes being removed are automatically relocated to other volumes within the policy.

The following example redefines a policy that has four volumes by adding two new volumes, removing an existing volume, and enforcing the policy for rebalancing.

To rebalance extents

- 1 Define the policy by specifying the `-o balance` and `-c` options:

```
# fsapadm define -o balance -c 2m /mnt loadbal vol1 vol2 vol4 \
vol5 vol6
```

- 2 Enforce the policy:

```
# fsapadm enforcefile -f strict /mnt/fileddb
```

Converting a multi-volume file system to a single volume file system

Because data can be relocated among volumes in a multi-volume file system, you can convert a multi-volume file system to a traditional, single volume file system by moving all file system data onto a single volume. Such a conversion is useful to users who would like to try using a multi-volume file system or Dynamic Storage Tiering, but are not committed to using a multi-volume file system permanently.

See [“About Dynamic Storage Tiering”](#) on page 143.

There are three restrictions to this operation:

- The single volume must be the first volume in the volume set
- The first volume must have sufficient space to hold all of the data and file system metadata
- The volume cannot have any allocation policies that restrict the movement of data

Converting to a single volume file system

The following procedure converts an existing multi-volume file system, `/mnt1`, of the volume set `vset1`, to a single volume file system, `/mnt1`, on volume `vol1` in diskgroup `dg1`.

Note: Steps 5, 6, 7, and 8 are optional, and can be performed if you prefer to remove the wrapper of the volume set object.

Converting to a single volume file system

- 1 Determine if the first volume in the volume set, which is identified as device number 0, has the capacity to receive the data from the other volumes that will be removed:

```
# df /mnt1
/mnt1  (/dev/vx/dsk/dg1/vol1):16777216 blocks  3443528 files
```

- 2 If the first volume does not have sufficient capacity, grow the volume to a sufficient size:

```
# fsvoladm resize /mnt1 vol1 150g
```

- 3 Remove all existing allocation policies:

```
# fsppadm unassign /mnt1
```

- 4 Remove all volumes except the first volume in the volume set:

```
# fsvoladm remove /mnt1 vol2
# vxvset -g dg1 rmvol vset1 vol2
# fsvoladm remove /mnt1 vol3
# vxvset -g dg1 rmvol vset1 vol3
```

Before removing a volume, the file system attempts to relocate the files on that volume. Successful relocation requires space on another volume, and no allocation policies can be enforced that pin files to that volume. The time for the command to complete is proportional to the amount of data that must be relocated.

- 5 Unmount the file system:

```
# umount /mnt1
```

- 6 Remove the volume from the volume set:

```
# vxvset -g dg1 rmvol vset1 vol1
```

7 Edit the `/etc/fstab` file to replace the volume set name, `vset1`, with the volume device name, `vol1`.

8 Mount the file system:

```
# mount -F vxfs /dev/vx/dsk/dg1/vol1 /mnt1
```


Dynamic Storage Tiering

This chapter includes the following topics:

- [About Dynamic Storage Tiering](#)
- [Placement classes](#)
- [Administering placement policies](#)
- [File placement policy grammar](#)
- [File placement policy rules](#)
- [Calculating I/O temperature and access temperature](#)
- [Multiple criteria in file placement policy rule statements](#)
- [File placement policy rule and statement ordering](#)
- [File placement policies and extending files](#)

About Dynamic Storage Tiering

VxFS uses multi-tier online storage by way of the Dynamic Storage Tiering (DST) feature, which functions on top of multi-volume file systems. Multi-volume file systems are file systems that occupy two or more virtual volumes. The collection of volumes is known as a volume set. A volume set is made up of disks or disk array LUNs belonging to a single Veritas Volume Manager (VxVM) disk group. A multi-volume file system presents a single name space, making the existence of multiple volumes transparent to users and applications. Each volume retains a separate identity for administrative purposes, making it possible to control the locations to which individual files are directed.

See [“About multi-volume support”](#) on page 122.

Note: Some of the commands have changed or been removed between the 4.1 release and the 5.0 release to make placement policy management more user-friendly. The following commands have been removed: `fsrpadm`, `fsmove`, and `fssweep`. The output of the `queryfile`, `queryfs`, and `list` options of the `fsapadm` command now print the allocation order by name instead of number.

DST allows administrators of multi-volume VxFS file systems to manage the placement of files on individual volumes in a volume set by defining placement policies. Placement policies control both initial file location and the circumstances under which existing files are relocated. These placement policies cause the files to which they apply to be created and extended on specific subsets of a file system's volume set, known as placement classes. The files are relocated to volumes in other placement classes when they meet the specified naming, timing, access rate, and storage capacity-related conditions.

You make a VxVM volume part of a placement class by associating a volume tag with it. For file placement purposes, VxFS treats all of the volumes in a placement class as equivalent, and balances space allocation across them. A volume may have more than one tag associated with it. If a volume has multiple tags, the volume belongs to multiple placement classes and is subject to allocation and relocation policies that relate to any of the placement classes. Multiple tagging should be used carefully.

See “Placement classes” on page 145.

VxFS imposes no capacity, performance, availability, or other constraints on placement classes. Any volume may be added to any placement class, no matter what type the volume has nor what types other volumes in the class have. However, a good practice is to place volumes of similar I/O performance and availability in the same placement class.

Note: Dynamic Storage Tiering is a licensed feature. You must purchase a separate license key for DST to operate. See the *Veritas Storage Foundation Release Notes*.

The *Using Dynamic Storage Tiering* Symantec Yellow Book provides additional information regarding the Dynamic Storage Tiering feature, including the value of DST and best practices for using DST. You can download *Using Dynamic Storage Tiering* from the following Web page:

<http://www.symantec.com/enterprise/yellowbooks/index.jsp>

Placement classes

A placement class is a Dynamic Storage Tiering attribute of a given volume in a volume set of a multi-volume file system. This attribute is a character string, and is known as a volume tag. A volume can have different tags, one of which can be the placement class. The placement class tag makes a volume distinguishable by DST.

Volume tags are organized as hierarchical name spaces in which periods separate the levels of the hierarchy. By convention, the uppermost level in the volume tag hierarchy denotes the Storage Foundation component or application that uses a tag, and the second level denotes the tag's purpose. DST recognizes volume tags of the form `vxfs.placement_class.class_name`. The prefix `vxfs` identifies a tag as being associated with VxFS. `placement_class` identifies the tag as a file placement class that DST uses. `class_name` represents the name of the file placement class to which the tagged volume belongs. For example, a volume with the tag `vxfs.placement_class.tier1` belongs to placement class `tier1`.

Administrators use the `vxassist` command to associate tags with volumes.

See the `vxassist(1M)` manual page.

VxFS policy rules specify file placement in terms of placement classes rather than in terms of individual volumes. All volumes that belong to a particular placement class are interchangeable with respect to file creation and relocation operations. Specifying file placement in terms of placement classes rather than in terms of specific volumes simplifies the administration of multi-tier storage.

The administration of multi-tier storage is simplified in the following ways:

- Adding or removing volumes does not require a file placement policy change. If a volume with a tag value of `vxfs.placement_class.tier2` is added to a file system's volume set, all policies that refer to `tier2` immediately apply to the newly added volume with no administrative action. Similarly, volumes can be evacuated, that is, have data removed from them, and be removed from a file system without a policy change. The active policy continues to apply to the file system's remaining volumes.
- File placement policies are not specific to individual file systems. A file placement policy can be assigned to any file system whose volume set includes volumes tagged with the tag values (placement classes) named in the policy. This property makes it possible for data centers with large numbers of servers to define standard placement policies and apply them uniformly to all servers with a single administrative action.

Tagging volumes as placement classes

The following example tags the `vsavola` volume as placement class `tier1`, `vsavolb` as placement class `tier2`, `vsavolc` as placement class `tier3`, and `vsavold` as placement class `tier4` using the `vxadm` command.

To tag volumes

- ◆ Tag the volumes as placement classes:

```
# vxassist -g cfsdg settag vsavola vxfs.placement_class.tier1
# vxassist -g cfsdg settag vsavolb vxfs.placement_class.tier2
# vxassist -g cfsdg settag vsavolc vxfs.placement_class.tier3
# vxassist -g cfsdg settag vsavold vxfs.placement_class.tier4
```

Listing placement classes

Placement classes are listed using the `vxassist listtagvxassist listtag` command.

See the `vxassist(1M)` manual page.

The following example lists all volume tags, including placement classes, set on a volume `vsavola` in the diskgroup `cfsdg`.

To list placement classes

- ◆ List the volume tags, including placement classes:

```
# vxassist -g cfsdg listtag vsavola
```

Administering placement policies

A VxFS file placement policy document contains rules by which VxFS creates, relocates, and deletes files, but the placement policy does not refer to specific file systems or volumes. You can create a file system's active file placement policy by assigning a placement policy document to the file system via the `fspadm` command or the GUI.

See the `fspadm(1M)` manual page.

At most, one file placement policy can be assigned to a VxFS file system at any time. A file system may have no file placement policy assigned to it, in which case VxFS allocates space for new files according to its own internal algorithms.

In systems with Storage Foundation Management Server (SFMS) software installed, file placement policy information is stored in the SFMS database. The SFMS

database contains both XML policy documents and lists of hosts and file systems for which each document is the current active policy. When a policy document is updated, SFMS can assign the updated document to all file systems whose current active policies are based on that document. By default, SFMS does not update file system active policies that have been created or modified locally, that is by the hosts that control the placement policies' file systems. If a SFMS administrator forces assignment of a placement policy to a file system, the file system's active placement policy is overwritten and any local changes that had been made to the placement policy are lost.

Assigning a placement policy

The following example uses the `fsppadm assign` command to assign the file placement policy represented in the XML policy document `/tmp/policy1.xml` for the file system at mount point `/mnt1`.

To assign a placement policy

- ◆ Assign a placement policy to a file system:

```
# fsppadm assign /mnt1 /tmp/policy1.xml
```

Unassigning a placement policy

The following example uses the `fsppadm unassign` command to unassign the active file placement policy from the file system at mount point `/mnt1`.

To unassign a placement policy

- ◆ Unassign the placement policy from a file system:

```
# fsppadm unassign /mnt1
```

Analyzing the space impact of enforcing a placement policy

The following example uses the `fsppadm analyze` command to analyze the impact if the enforce operation is performed on the file placement policy represented in the XML policy document `/tmp/policy1.xml` for the mount point `/mnt1`. The command builds the I/O temperature database if necessary.

To analyze the space impact of enforcing a placement policy

- ◆ Analyze the impact of enforcing the file placement policy represented in the XML policy document `/tmp/policy1.xml` for the mount point `/mnt1`:

```
# fsppadm analyze -F /tmp/policy1.xml -i /mnt1
```

Querying which files will be affected by enforcing a placement policy

The following example uses the `fsppadm query` command to generate a list of files that will be affected by enforcing a placement policy. The command provides details about where the files currently reside, to where the files will be relocated, and which rule in the placement policy applies to the files.

To query which files will be affected by enforcing a placement policy

- ◆ Query the files that will be affected:

```
# fsppadm query /mnt1/dir1/dir2 /mnt2 /mnt1/dir3
```

Enforcing a placement policy

Enforcing a placement policy for a file system requires that the policy be assigned to the file system. You must assign a placement policy before it can be enforced.

See [“Assigning a placement policy”](#) on page 147.

Enforce operations are logged in a hidden file, `.__fsppadm_enforce.log`, in the `lost+found` directory of the mount point. This log file contains details such as files' previous locations, the files' new locations, and the reasons for the files' relocations. The enforce operation creates the `.__fsppadm_enforce.log` file if the file does not exist. The enforce operation appends the file if the file already exists. The `.__fsppadm_enforce.log` file can be backed up or removed as with a normal file.

You can specify the `-F` option to specify a placement policy other than the existing active placement policy. This option can be used to enforce the rules given in the specified placement policy for maintenance purposes, such as for reclaiming a LUN from the file system.

You can specify the `-p` option to specify the number of concurrent threads to be used to perform the `fsppadm` operation. You specify the `io_nice` parameter as an integer between 1 and 100, with 50 being the default value. A value of 1 specifies 1 slave and 1 master thread per mount. A value of 50 specifies 16 slaves and 1 master thread per mount. A value of 100 specifies 32 slaves and 1 master thread per mount.

You can specify the `-c` option so that the `fsppadm` command processes only those files that have some activity stats logged in the File Change Log (FCL) file during the period specified in the placement policy. You can use the `-c` option only if the policy's `ACCESSTEMP` or `IOTEMP` elements use the `Prefer` criteria.

You can specify the `-T` option to specify the placement classes that contain files for the `fsppadm` command to sweep and relocate selectively. You can specify the `-T` option only if the policy uses the `Prefer` criteria for `IOTEMP`.

See the `fsppadm(1M)` manual page.

The following example uses the `fsppadm enforce` command to enforce the file placement policy for the file system at mount point `/mnt1`, and includes the access time, modification time, and file size of the specified paths in the report, `/tmp/report`.

To enforce a placement policy

- ◆ Enforce a placement policy for a file system:

```
# fsppadm enforce -a -r /tmp/report /mnt1
Current Current Relocated Relocated
Class Volume Class Volume Rule File
tier3 dstvole tier3 dstvole a_to_z /mnt1/mds1/d1/file1
tier3 dstvole tier3 dstvole a_to_z /mnt1/mds1/d1/file2
tier3 dstvole tier3 dstvole a_to_z /mnt1/mds1/d1/d2/file3
tier3 dstvolf tier3 dstvolf a_to_z /mnt1/mds1/d1/d2/file4
.
.
.
Sweep path : /mnt1
Files moved : 42
KB moved : 1267
```

Tier Name	Size (KB)	Free Before (KB)	Free After (KB)
tier4	524288	524256	524256
tier3	524288	522968	522968
tier2	524288	524256	524256
tier1	524288	502188	501227

Validating a placement policy

The following example uses the `fsppadm validate` command to validate the placement policy `policy.xml` against all mounted file systems.

To validate a placement policy against all mounted file systems

- ◆ Validate the placement policy:

```
# fsppadm validate /tmp/policy.xml
```

File placement policy grammar

VxFS allocates and relocates files within a multi-volume file system based on properties in the file system metadata that pertains to the files. Placement decisions may be based on file name, directory of residence, time of last access, access frequency, file size, and ownership. An individual file system's criteria for allocating and relocating files are expressed in the file system's file placement policy.

A VxFS file placement policy defines the desired placement of sets of files on the volumes of a VxFS multi-volume file system. A file placement policy specifies the placement classes of volumes on which files should be created, and where and under what conditions the files should be relocated to volumes in alternate placement classes or deleted. You can create file placement policy documents, which are XML text files, using either an XML or text editor, or a VxFS graphical interface wizard.

See the `/opt/VRTSvxfs/etc/placement_policy.dtd` file for the overall structure of a placement policy.

File placement policy rules

A VxFS file placement policy consists of one or more rules. Each rule applies to one or more files. The files to which a rule applies are designated in one or more `SELECT` statements. A `SELECT` statement designates files according to one or more of four properties: their names or naming patterns, the directories in which they reside, their owners' user names, and their owners' group names.

A file may be designated by more than one rule. For example, if one rule designates files in directory `/dir`, and another designates files owned by `user1`, a file in `/dir` that is owned by `user1` is designated by both rules. Only the rule that appears first in the placement policy applies to the file; subsequent rules are ignored.

You can define placement policies that do not encompass the entire file system name space. When a file that is not designated by any rule in its file system's active placement policy is created, VxFS places the file according to its own internal algorithms. To maintain full control over file placement, include a catchall rule at the end of each placement policy document with a `SELECT` statement that designates files by the naming pattern `*`. Such a rule designates all files that have not been designated by the rules appearing earlier in the placement policy document.

Two types of rules exist: `data` and `ckpt`. The `data` rule type allows DST to relocate normal data files. The `ckpt` rule type allows DST to relocate Storage Checkpoints. You specify the rule type by setting the `Flags` attribute for the rule.

SELECT statement

The VxFS placement policy rule `SELECT` statement designates the collection of files to which a rule applies.

The following XML snippet illustrates the general form of the `SELECT` statement:

```
<SELECT>
  <DIRECTORY Flags="directory_flag_value"> value
</DIRECTORY>
  <PATTERN Flags="pattern_flag_value"> value </PATTERN>
  <USER> value </USER>
  <GROUP> value </GROUP>
</SELECT>
```

A `SELECT` statement may designate files by using the following selection criteria:

`<DIRECTORY>` A full path name relative to the file system mount point. The `Flags="directory_flag_value"` XML attribute must have a value of `nonrecursive`, denoting that only files in the specified directory are designated, or a value of `recursive`, denoting that files in all subdirectories of the specified directory are designated. The `Flags` attribute is mandatory.

The `<DIRECTORY>` criterion is optional, and may be specified more than once.

<PATTERN>	<p>Either an exact file name or a pattern using a single wildcard character (*). For example, the pattern "abc*" denotes all files whose names begin with "abc". The pattern "abc.*" denotes all files whose names are exactly "abc" followed by a period and any extension. The pattern "**abc" denotes all files whose names end in "abc", even if the name is all or part of an extension. The pattern "*.abc" denotes files of any name whose name extension (following the period) is "abc". The pattern "ab*c" denotes all files whose names start with "ab" and end with "c". The first "*" character is treated as a wildcard, while any subsequent "*" characters are treated as literal text. The pattern cannot contain "/".</p> <p>The wildcard character matches any character, including ".", "?", and "[", unlike using the wildcard in a shell.</p> <p>The <code>Flags="pattern_flag_value"</code> XML attribute is optional, and if specified can only have a value of <code>recursive</code>. Specify <code>Flags="recursive"</code> only if the pattern is a directory. If <code>Flags</code> is not specified, the default attribute value is <code>nonrecursive</code>. If <code>Flags="recursive"</code> is specified, the enclosing selection criteria selects all files in any component directory that is anywhere below the directory specified by <DIRECTORY> if the component directory matches the pattern and either of the following is true:</p> <ul style="list-style-type: none"> ■ <DIRECTORY> is specified and has the recursive flag. ■ <DIRECTORY> is not specified and the directory is anywhere in the file system. <p>If the pattern contains the wildcard character (*), wildcard character matching is performed.</p> <p>The <PATTERN> criterion is optional, and may be specified more than once. Only one value can be specified per <PATTERN> element.</p>
<USER>	<p>User name of the file's owner. The user number cannot be specified in place of the name.</p> <p>The <USER> criterion is optional, and may be specified more than once.</p>
<GROUP>	<p>Group name of the file's owner. The group number cannot be specified in place of the group name.</p> <p>The <GROUP> criterion is optional, and may be specified more than once.</p>

One or more instances of any or all of the file selection criteria may be specified within a single `SELECT` statement. If two or more selection criteria of different types are specified in a single statement, a file must satisfy one criterion of each type to be selected.

In the following example, only files that reside in either the `ora/db` or the `crash/dump` directory, and whose owner is either `user1` or `user2` are selected for possible action:

```
<SELECT>
  <DIRECTORY Flags="nonrecursive">ora/db</DIRECTORY>
  <DIRECTORY Flags="nonrecursive">crash/dump</DIRECTORY>
  <USER>user1</USER>
  <USER>user2</USER>
</SELECT>
```

A rule may include multiple `SELECT` statements. If a file satisfies the selection criteria of one of the `SELECT` statements, it is eligible for action.

In the following example, any files owned by either `user1` or `user2`, no matter in which directories they reside, as well as all files in the `ora/db` or `crash/dump` directories, no matter which users own them, are eligible for action:

```
<SELECT>
  <DIRECTORY Flags="nonrecursive">ora/db</DIRECTORY>
  <DIRECTORY Flags="nonrecursive">crash/dump</DIRECTORY>
</SELECT>
<SELECT>
  <USER>user1</USER>
  <USER>user2</USER>
</SELECT>
```

When VxFS creates new files, VxFS applies active placement policy rules in the order of appearance in the active placement policy's XML source file. The first rule in which a `SELECT` statement designates the file to be created determines the file's placement; no later rules apply. Similarly, VxFS scans the active policy rules on behalf of each file when relocating files, stopping the rules scan when it reaches the first rule containing a `SELECT` statement that designates the file. This behavior holds true even if the applicable rule results in no action. Take for example a policy rule that indicates that `.dat` files inactive for 30 days should be relocated, and a later rule indicates that `.dat` files larger than 10 megabytes should be relocated. A 20 megabyte `.dat` file that has been inactive for 10 days will not be relocated because the earlier rule applied. The later rule is never scanned.

A placement policy rule's action statements apply to all files designated by any of the rule's `SELECT` statements. If an existing file is not designated by a `SELECT` statement in any rule of a file system's active placement policy, then DST does not relocate or delete the file. If an application creates a file that is not designated by a `SELECT` statement in a rule of the file system's active policy, then VxFS places the file according to its own internal algorithms. If this behavior is inappropriate,

the last rule in the policy document on which the file system's active placement policy is based should specify `<PATTERN>*``</PATTERN>` as the only selection criterion in its `SELECT` statement, and a `CREATE` statement naming the desired placement class for files not selected by other rules.

CREATE statement

A `CREATE` statement in a file placement policy rule specifies one or more placement classes of volumes on which VxFS should allocate space for new files to which the rule applies at the time the files are created. You can specify only placement classes, not individual volume names, in a `CREATE` statement.

A file placement policy rule may contain at most one `CREATE` statement. If a rule does not contain a `CREATE` statement, VxFS places files designated by the rule's `SELECT` statements according to its internal algorithms. However, rules without `CREATE` statements can be used to relocate or delete existing files that the rules' `SELECT` statements designate.

The following XML snippet illustrates the general form of the `CREATE` statement:

```
<CREATE>
  <ON Flags="flag_value">
    <DESTINATION>
      <CLASS> placement_class_name </CLASS>
      <BALANCE_SIZE Units="units_specifier"> chunk_size
    </BALANCE_SIZE>
    </DESTINATION>
    <DESTINATION> additional_placement_class_specifications
  </DESTINATION>
  </ON>
</CREATE>
```

A `CREATE` statement includes a single `<ON>` clause, in which one or more `<DESTINATION>` XML elements specify placement classes for initial file allocation in order of decreasing preference. VxFS allocates space for new files to which a rule applies on a volume in the first class specified, if available space permits. If space cannot be allocated on any volume in the first class, VxFS allocates space on a volume in the second class specified if available space permits, and so forth.

If space cannot be allocated on any volume in any of the placement classes specified, file creation fails with an `ENOSPC` error, even if adequate space is available elsewhere in the file system's volume set. This situation can be circumvented by specifying a `Flags` attribute with a value of "any" in the `<ON>` clause. If `<ON Flags="any">` is specified in a `CREATE` statement, VxFS first attempts to allocate

space for new files to which the rule applies on the specified placement classes. Failing that, VxFS resorts to its internal space allocation algorithms, so file allocation does not fail unless there is no available space any-where in the file system's volume set.

The `Flags="any"` attribute differs from the catchall rule in that this attribute applies only to files designated by the `SELECT` statement in the rule, which may be less inclusive than the `<PATTERN>*</PATTERN>` file selection specification of the catchall rule.

In addition to the placement class name specified in the `<CLASS>` sub-element, a `<DESTINATION>` XML element may contain a `<BALANCE_SIZE>` sub-element. Presence of a `<BALANCE_SIZE>` element indicates that space allocation should be distributed across the volumes of the placement class in chunks of the indicated size. For example, if a balance size of one megabyte is specified for a placement class containing three volumes, VxFS allocates the first megabyte of space for a new or extending file on the first (lowest indexed) volume in the class, the second megabyte on the second volume, the third megabyte on the third volume, the fourth megabyte on the first volume, and so forth. Using the `Units` attribute in the `<BALANCE_SIZE>` XML tag, the balance size value may be specified in the following units:

bytes	Bytes
KB	Kilobytes
MB	Megabytes
GB	Gigabytes

The `<BALANCE_SIZE>` element distributes the allocation of database files across the volumes in a placement class. In principle, distributing the data in each file across multiple volumes distributes the I/O load across the volumes as well.

The `CREATE` statement in the following example specifies that files to which the rule applies should be created on the `tier1` volume if space is available, and on one of the `tier2` volumes if not. If space allocation on `tier1` and `tier2` volumes is not possible, file creation fails, even if space is available on `tier3` volumes.

```
<CREATE>
  <ON>
    <DESTINATION>
      <CLASS>tier1</CLASS>
    </DESTINATION>
  <DESTINATION>
```

```

        <CLASS>tier2</CLASS>
        <BALANCE_SIZE Units="MB">1</BALANCE_SIZE>
    </DESTINATION>
</ON>
</CREATE>

```

The `<BALANCE_SIZE>` element with a value of one megabyte is specified for allocations on `tier2` volumes. For files allocated on `tier2` volumes, the first megabyte would be allocated on the first volume, the second on the second volume, and so forth.

RELOCATE statement

The `RELOCATE` action statement of file placement policy rules specifies an action that VxFS takes on designated files during periodic scans of the file system, and the circumstances under which the actions should be taken. The `fspadm enforce` command is used to scan all or part of a file system for files that should be relocated based on rules in the active placement policy at the time of the scan.

See the `fspadm(1M)` manual page.

The `fspadm enforce` command scans file systems in path name order. For each file, VxFS identifies the first applicable rule in the active placement policy, as determined by the rules' `SELECT` statements. If the file resides on a volume specified in the `<FROM>` clause of one of the rule's `RELOCATE` statements, and if the file meets the criteria for relocation specified in the statement's `<WHEN>` clause, the file is scheduled for relocation to a volume in the first placement class listed in the `<TO>` clause that has space available for the file. The scan that results from issuing the `fspadm enforce` command runs to completion before any files are relocated.

The following XML snippet illustrates the general form of the `RELOCATE` statement:

```

<RELOCATE>
  <FROM>
    <SOURCE>
      <CLASS> placement_class_name </CLASS>
    </SOURCE>
    <SOURCE> additional_placement_class_specifications
    </SOURCE>
  </FROM>
  <TO>
    <DESTINATION>
      <CLASS> placement_class_name </CLASS>
      <BALANCE_SIZE Units="units_specifier">
        chunk_size

```

```
        </BALANCE_SIZE>
    </DESTINATION>
    <DESTINATION>
        additional_placement_class_specifications
    </DESTINATION>
</TO>
<WHEN> relocation_conditions </WHEN>
</RELOCATE>
```

A **RELOCATE** statement contains the following clauses:

<FROM>	<p>An optional clause that contains a list of placement classes from whose volumes designated files should be relocated if the files meet the conditions specified in the <WHEN> clause. No priority is associated with the ordering of placement classes listed in a <FROM> clause. If a file to which the rule applies is located on a volume in any specified placement class, the file is considered for relocation.</p> <p>If a RELOCATE statement contains a <FROM> clause, VxFS only considers files that reside on volumes in placement classes specified in the clause for relocation. If no <FROM> clause is present, qualifying files are relocated regardless of where the files reside.</p>
---------------------	---

<TO>

Indicates the placement classes to which qualifying files should be relocated. Unlike the source placement class list in a FROM clause, placement classes in a <TO> clause are specified in priority order. Files are relocated to volumes in the first specified placement class if possible, to the second if not, and so forth.

The <TO> clause of the RELOCATE statement contains a list of <DESTINATION> XML elements specifying placement classes to whose volumes VxFS relocates qualifying files. Placement classes are specified in priority order. VxFS relocates qualifying files to volumes in the first placement class specified as long as space is available. A <DESTINATION> element may contain an optional <BALANCE_SIZE> modifier sub-element. The <BALANCE_SIZE> modifier indicates that relocated files should be distributed across the volumes of the destination placement class in chunks of the indicated size. For example, if a balance size of one megabyte is specified for a placement class containing three volumes, VxFS relocates the first megabyte the file to the first (lowest indexed) volume in the class, the second megabyte to the second volume, the third megabyte to the third volume, the fourth megabyte to the first volume, and so forth. Using the Units attribute in the <BALANCE_SIZE> XML tag, the chunk value may be specified in the balance size value may be specified in bytes (Units="bytes"), kilobytes (Units="KB"), megabytes (Units="MB"), or gigabytes (Units="GB").

The <BALANCE_SIZE> element distributes the allocation of database files across the volumes in a placement class. In principle, distributing the data in each file across multiple volumes distributes the I/O load across the volumes as well.

<WHEN>

An optional clause that indicates the conditions under which files to which the rule applies should be relocated. Files that have been unaccessed or unmodified for a specified period, reached a certain size, or reached a specific I/O temperature or access temperature level may be relocated. If a RELOCATE statement does not contain a <WHEN> clause, files to which the rule applies are relocated unconditionally.

A <WHEN> clause may be included in a RELOCATE statement to specify that files should be relocated only if any or all of four types of criteria are met. Files can be specified for relocation if they satisfy one or more criteria.

The following are the criteria that can be specified for the <WHEN> clause:

<ACCAGE>

This criterion is met when files are inactive for a designated period or during a designated period relative to the time at which the `fspadm enforce` command was issued.

<MODAGE>	This criterion is met when files are unmodified for a designated period or during a designated period relative to the time at which the <code>fsppadm enforce</code> command was issued.
<SIZE>	This criterion is met when files exceed or drop below a designated size or fall within a designated size range.
<IOTEMP>	<p>This criterion is met when files exceed or drop below a designated I/O temperature, or fall within a designated I/O temperature range. A file's I/O temperature is a measure of the I/O activity against it during the period designated by the <PERIOD> element prior to the time at which the <code>fsppadm enforce</code> command was issued.</p> <p>See “Calculating I/O temperature and access temperature” on page 169.</p>
<ACCESSTEMP>	This criterion is met when files exceed or drop below a specified average access temperature, or fall within a specified access temperature range. A file's access temperature is similar to its I/O temperature, except that access temperature is computed using the number of I/O requests to the file, rather than the number of bytes transferred.

Note: The use of <IOTEMP> and <ACCESSTEMP> for data placement on VxFS servers that are used as NFS servers may not be very effective due to NFS caching. NFS client side caching and the way that NFS works can result in I/O initiated from an NFS client not producing NFS server side I/O. As such, any temperature measurements in place on the server side will not correctly reflect the I/O behavior that is specified by the placement policy.

If the server is solely used as an NFS server, this problem can potentially be mitigated by suitably adjusting or lowering the temperature thresholds. However, adjusting the thresholds may not always create the desired effect. In addition, if the same mount point is used both as an NFS export as well as a local mount, the temperature-based placement decisions will not be very effective due to the NFS cache skew.

The following XML snippet illustrates the general form of the <WHEN> clause in a `RELOCATE` statement:

```
<WHEN>
  <ACCAGE Units="units_value">
    <MIN Flags="comparison_operator">
      min_access_age</MIN>
    <MAX Flags="comparison_operator">
```

```

        max_access_age</MAX>
    </ACCAGE>
    <MODAGE Units="units_value">
        <MIN Flags="comparison_operator">
            min_modification_age</MIN>
        <MAX Flags="comparison_operator">
            max_modification_age</MAX>
    </MODAGE>
    <SIZE " Units="units_value">
        <MIN Flags="comparison_operator">
            min_size</MIN>
        <MAX Flags="comparison_operator">
            max_size</MAX>
    </SIZE>
    <IOTEMP Type="read_write_preference" Prefer="temperature_preference">
        <MIN Flags="comparison_operator">
            min_I/O_temperature</MIN>
        <MAX Flags="comparison_operator">
            max_I/O_temperature</MAX>
        <PERIOD Units="days_or_hours"> days_or_hours_of_interest </PERIOD>
    </IOTEMP>
    <ACCESSTEMP Type="read_write_preference"
        Prefer="temperature_preference">
        <MIN Flags="comparison_operator">
            min_access_temperature</MIN>
        <MAX Flags="comparison_operator">
            max_access_temperature</MAX>
        <PERIOD Units="days_or_hours"> days_or_hours_of_interest </PERIOD>
    </ACCESSTEMP>
</WHEN>

```

The access age (<ACCAGE>) element refers to the amount of time since a file was last accessed. VxFS computes access age by subtracting a file's time of last access, `atime`, from the time when the `fsppadm enforce` command was issued. The <MIN> and <MAX> XML elements in an <ACCAGE> clause, denote the minimum and maximum access age thresholds for relocation, respectively. These elements are optional, but at least one must be included. Using the `Units` XML attribute, the <MIN> and <MAX> elements may be specified in the following units:

hours	Hours
days	Days. A day is considered to be 24 hours prior to the time that the <code>fsppadm enforce</code> command was issued.

Both the `<MIN>` and `<MAX>` elements require `Flags` attributes to direct their operation.

For `<MIN>`, the following `Flags` attributes values may be specified:

<code>gt</code>	The time of last access must be greater than the specified interval.
<code>eq</code>	The time of last access must be equal to the specified interval.
<code>gteq</code>	The time of last access must be greater than or equal to the specified interval.

For `<MAX>`, the following `Flags` attributes values may be specified.

<code>lt</code>	The time of last access must be less than the specified interval.
<code>lteq</code>	The time of last access must be less than or equal to the specified interval.

Including a `<MIN>` element in a `<WHEN>` clause causes VxFS to relocate files to which the rule applies that have been inactive for longer than the specified interval. Such a rule would typically be used to relocate inactive files to less expensive storage tiers. Conversely, including `<MAX>` causes files accessed within the specified interval to be relocated. It would typically be used to move inactive files against which activity had recommenced to higher performance or more reliable storage. Including both `<MIN>` and `<MAX>` causes VxFS to relocate files whose access age lies between the two.

The modification age relocation criterion, `<MODAGE>`, is similar to access age, except that files' POSIX `mtime` values are used in computations. You would typically specify the `<MODAGE>` criterion to cause relocation of recently modified files to higher performance or more reliable storage tiers in anticipation that the files would be accessed recurrently in the near future.

The file size relocation criterion, `<SIZE>`, causes files to be relocated if the files are larger or smaller than the values specified in the `<MIN>` and `<MAX>` relocation criteria, respectively, at the time that the `fsppadm enforce` command was issued. Specifying both criteria causes VxFS to schedule relocation for files whose sizes lie between the two. Using the `Units` attribute, threshold file sizes may be specified in the following units:

<code>bytes</code>	Bytes
<code>KB</code>	Kilobytes
<code>MB</code>	Megabytes

GB

Gigabytes

Specifying the I/O temperature relocation criterion

The I/O temperature relocation criterion, `<IOTEMP>`, causes files to be relocated if their I/O temperatures rise above or drop below specified values over a specified period immediately prior to the time at which the `fsppadm enforce` command was issued. A file's I/O temperature is a measure of the read, write, or total I/O activity against it normalized to the file's size. Higher I/O temperatures indicate higher levels of application activity; lower temperatures indicate lower levels. VxFS computes a file's I/O temperature by dividing the number of bytes transferred to or from it (read, written, or both) during the specified period by its size at the time that the `fsppadm enforce` command was issued.

See [“Calculating I/O temperature and access temperature”](#) on page 169.

As with the other file relocation criteria, `<IOTEMP>` may be specified with a lower threshold by using the `<MIN>` element, an upper threshold by using the `<MAX>` element, or as a range by using both. However, I/O temperature is dimensionless and therefore has no specification for units.

VxFS computes files' I/O temperatures over the period between the time when the `fsppadm enforce` command was issued and the number of days or hours in the past specified in the `<PERIOD>` element, where a day is a 24 hour period. The default unit of time is days. You can specify hours as the time unit by setting the `Units` attribute of the `<PERIOD>` element to `hours`. Symantec recommends that you specify hours only if you are using solid state disks (SSDs).

For example, if you issued the `fsppadm enforce` command at 2 PM on Wednesday and you want VxFS to look at file I/O activity for the period between 2 PM on Monday and 2 PM on Wednesday, which is a period of 2 days, you would specify the following `<PERIOD>` element:

```
<PERIOD> 2 </PERIOD>
```

If you instead want VxFS to look at file I/O activity between 3 hours prior to running the `fsppadm enforce` command and the time that you ran the command, you specify the following `<PERIOD>` element:

```
<PERIOD Units="hours"> 3 </PERIOD>
```

The amount of time specified in the `<PERIOD>` element should not exceed one or two weeks due to the disk space used by the File Change Log (FCL) file.

See [“About the File Change Log file”](#) on page 114.

I/O temperature is a softer measure of I/O activity than access age. With access age, a single access to a file resets the file's atime to the current time. In contrast, a file's I/O temperature decreases gradually as time passes without the file being accessed, and increases gradually as the file is accessed periodically. For example, if a new 10 megabyte file is read completely five times on Monday and `fsppadm enforce` runs at midnight, the file's two-day I/O temperature will be five and its access age in days will be zero. If the file is read once on Tuesday, the file's access age in days at midnight will be zero, and its two-day I/O temperature will have dropped to three. If the file is read once on Wednesday, the file's access age at midnight will still be zero, but its two-day I/O temperature will have dropped to one, as the influence of Monday's I/O will have disappeared.

If the intention of a file placement policy is to keep files in place, such as on top-tier storage devices, as long as the files are being accessed at all, then access age is the more appropriate relocation criterion. However, if the intention is to relocate files as the I/O load on them decreases, then I/O temperature is more appropriate.

The case for upward relocation is similar. If files that have been relocated to lower-tier storage devices due to infrequent access experience renewed application activity, then it may be appropriate to relocate those files to top-tier devices. A policy rule that uses access age with a low `<MAX>` value, that is, the interval between `fsppadm enforce` runs, as a relocation criterion will cause files to be relocated that have been accessed even once during the interval. Conversely, a policy that uses I/O temperature with a `<MIN>` value will only relocate files that have experienced a sustained level of activity over the period of interest.

RELOCATE statement examples

The following example illustrates an unconditional relocation statement, which is the simplest form of the `RELOCATE` policy rule statement:

```
<RELOCATE>
  <FROM>
    <SOURCE>
      <CLASS>tier1</CLASS>
    </SOURCE>
  </FROM>
  <TO>
    <DESTINATION>
      <CLASS>tier2</CLASS>
    </DESTINATION>
  </TO>
</RELOCATE>
```

The files designated by the rule's `SELECT` statement that reside on volumes in placement class `tier1` at the time the `fsppadm enforce` command executes would be unconditionally relocated to volumes in placement class `tier2` as long as space permitted. This type of rule might be used, for example, with applications that create and access new files but seldom access existing files once they have been processed. A `CREATE` statement would specify creation on `tier1` volumes, which are presumably high performance or high availability, or both. Each instantiation of `fsppadm enforce` would relocate files created since the last run to `tier2` volumes.

The following example illustrates a more comprehensive form of the `RELOCATE` statement that uses access age as the criterion for relocating files from `tier1` volumes to `tier2` volumes. This rule is designed to maintain free space on `tier1` volumes by relocating inactive files to `tier2` volumes:

```
<RELOCATE>
  <FROM>
    <SOURCE>
      <CLASS>tier1</CLASS>
    </SOURCE>
  </FROM>
  <TO>
    <DESTINATION>
      <CLASS>tier2</CLASS>
    </DESTINATION>
  </TO>
  <WHEN>
    <SIZE Units="MB">
      <MIN Flags="gt">1</MIN>
      <MAX Flags="lt">1000</MAX>
    </SIZE>
    <ACCAGE Units="days">
      <MIN Flags="gt">30</MIN>
    </ACCAGE>
  </WHEN>
</RELOCATE>
```

Files designated by the rule's `SELECT` statement are relocated from `tier1` volumes to `tier2` volumes if they are between 1 MB and 1000 MB in size and have not been accessed for 30 days. VxFS relocates qualifying files in the order in which it encounters them as it scans the file system's directory tree. VxFS stops scheduling qualifying files for relocation when it calculates that already-scheduled relocations would result in `tier2` volumes being fully occupied.

The following example illustrates a possible companion rule that relocates files from `tier2` volumes to `tier1` ones based on their I/O temperatures. This rule might be used to return files that had been relocated to `tier2` volumes due to inactivity to `tier1` volumes when application activity against them increases. Using I/O temperature rather than access age as the relocation criterion reduces the chance of relocating files that are not actually being used frequently by applications. This rule does not cause files to be relocated unless there is sustained activity against them over the most recent two-day period.

```
<RELOCATE>
  <FROM>
    <SOURCE>
      <CLASS>tier2</CLASS>
    </SOURCE>
  </FROM>
  <TO>
    <DESTINATION>
      <CLASS>tier1</CLASS>
    </DESTINATION>
  </TO>
  <WHEN>
    <IOTEMP Type="nrbytes">
      <MIN Flags="gt">5</MIN>
      <PERIOD>2</PERIOD>
    </IOTEMP>
  </WHEN>
</RELOCATE>
```

This rule relocates files that reside on `tier2` volumes to `tier1` volumes if their I/O temperatures are above 5 for the two day period immediately preceding the issuing of the `fsppadm enforce` command. VxFS relocates qualifying files in the order in which it encounters them during its file system directory tree scan. When `tier1` volumes are fully occupied, VxFS stops scheduling qualifying files for relocation.

VxFS file placement policies are able to control file placement across any number of placement classes. The following example illustrates a rule for relocating files with low I/O temperatures from `tier1` volumes to `tier2` volumes, and to `tier3` volumes when `tier2` volumes are fully occupied:

```
<RELOCATE>
  <FROM>
    <SOURCE>
      <CLASS>tier1</CLASS>
```

```

    </SOURCE>
  </FROM>
  <TO>
    <DESTINATION>
      <CLASS>tier2</CLASS>
    </DESTINATION>
    <DESTINATION>
      <CLASS>tier3</CLASS>
    </DESTINATION>
  </TO>
  <WHEN>
    <IOTEMP Type="nrbytes">
      <MAX Flags="lt">4</MAX>
      <PERIOD>3</PERIOD>
    </IOTEMP>
  </WHEN>
</RELOCATE>

```

This rule relocates files whose 3-day I/O temperatures are less than 4 and which reside on `tier1` volumes. When VxFS calculates that already-relocated files would result in `tier2` volumes being fully occupied, VxFS relocates qualifying files to `tier3` volumes instead. VxFS relocates qualifying files as it encounters them in its scan of the file system directory tree.

The `<FROM>` clause in the `RELOCATE` statement is optional. If the clause is not present, VxFS evaluates files designated by the rule's `SELECT` statement for relocation no matter which volumes they reside on when the `fsppadm enforce` command is issued. The following example illustrates a fragment of a policy rule that relocates files according to their sizes, no matter where they reside when the `fsppadm enforce` command is issued:

```

<RELOCATE>
  <TO>
    <DESTINATION>
      <CLASS>tier1</CLASS>
    </DESTINATION>
  </TO>
  <WHEN>
    <SIZE Units="MB">
      <MAX Flags="lt">10</MAX>
    </SIZE>
  </WHEN>
</RELOCATE>
<RELOCATE>

```

```

<TO>
  <DESTINATION>
    <CLASS>tier2</CLASS>
  </DESTINATION>
</TO>
<WHEN>
  <SIZE Units="MB">
    <MIN Flags="gteq">10</MIN>
    <MAX Flags="lt">100</MAX>
  </SIZE>
</WHEN>
</RELOCATE>
<RELOCATE>
  <TO>
    <DESTINATION>
      <CLASS>tier3</CLASS>
    </DESTINATION>
  </TO>
  <WHEN>
    <SIZE Units="MB">
      <MIN Flags="gteq">100</MIN>
    </SIZE>
  </WHEN>
</RELOCATE>

```

This rule relocates files smaller than 10 megabytes to `tier1` volumes, files between 10 and 100 megabytes to `tier2` volumes, and files larger than 100 megabytes to `tier3` volumes. VxFS relocates all qualifying files that do not already reside on volumes in their `DESTINATION` placement classes when the `fspadm enforce` command is issued.

DELETE statement

The `DELETE` file placement policy rule statement is very similar to the `RELOCATE` statement in both form and function, lacking only the `<TO>` clause. File placement policy-based deletion may be thought of as relocation with a fixed destination.

Note: Use `DELETE` statements with caution.

The following XML snippet illustrates the general form of the `DELETE` statement:

```

<DELETE>
  <FROM>

```

```

<SOURCE>
  <CLASS> placement_class_name </CLASS>
</SOURCE>
<SOURCE>
  additional_placement_class_specifications
</SOURCE>
</FROM>
<WHEN> relocation_conditions </WHEN>
</DELETE>

```

A **DELETE** statement contains the following clauses:

<FROM>	An optional clause that contains a list of placement classes from whose volumes designated files should be deleted if the files meet the conditions specified in the <WHEN> clause. No priority is associated with the ordering of placement classes in a <FROM> clause. If a file to which the rule applies is located on a volume in any specified placement class, the file is deleted. If a DELETE statement does not contain a <FROM> clause, VxFS deletes qualifying files no matter on which of a file system's volumes the files reside.
<WHEN>	An optional clause specifying the conditions under which files to which the rule applies should be deleted. The form of the <WHEN> clause in a DELETE statement is identical to that of the <WHEN> clause in a RELOCATE statement. If a DELETE statement does not contain a <WHEN> clause, files designated by the rule's SELECT statement, and the <FROM> clause if it is present, are deleted unconditionally.

DELETE statement examples

The following example illustrates the use of the **DELETE** statement:

```

<DELETE>
  <FROM>
    <SOURCE>
      <CLASS>tier3</CLASS>
    </SOURCE>
  </FROM>
</DELETE>
<DELETE>
  <FROM>
    <SOURCE>
      <CLASS>tier2</CLASS>
    </SOURCE>
  </FROM>
</DELETE>

```



```

</FROM>
<WHEN>
  <ACCAGE Units="days">
    <MIN Flags="gt">120</MIN>
  </ACCAGE>
</WHEN>
</DELETE>

```

The first `DELETE` statement unconditionally deletes files designated by the rule's `SELECT` statement that reside on `tier3` volumes when the `fspadm enforce` command is issued. The absence of a `<WHEN>` clause in the `DELETE` statement indicates that deletion of designated files is unconditional.

The second `DELETE` statement deletes files to which the rule applies that reside on `tier2` volumes when the `fspadm enforce` command is issued and that have not been accessed for the past 120 days.

Calculating I/O temperature and access temperature

An important application of VxFS Dynamic Storage Tiering is automating the relocation of inactive files to lower cost storage. If a file has not been accessed for the period of time specified in the `<ACCAGE>` element, a scan of the file system should schedule the file for relocation to a lower tier of storage. But, time since last access is inadequate as the only criterion for activity-based relocation.

Why time since last access is inadequate as the only criterion for activity-based relocation:

- Access age is a binary measure. The time since last access of a file is computed by subtracting the time at which the `fspadm enforce` command is issued from the POSIX `atime` in the file's metadata. If a file is opened the day before the `fspadm enforce` command, its time since last access is one day, even though it may have been inactive for the month preceding. If the intent of a policy rule is to relocate inactive files to lower tier volumes, it will perform badly against files that happen to be accessed, however casually, within the interval defined by the value of the `<ACCAGE>` parameter.
- Access age is a poor indicator of resumption of significant activity. Using `ACCAGE`, the time since last access, as a criterion for relocating inactive files to lower tier volumes may fail to schedule some relocations that should be performed, but at least this method results in less relocation activity than necessary. Using `ACCAGE` as a criterion for relocating previously inactive files that have become active is worse, because this method is likely to schedule relocation activity that is not warranted. If a policy rule's intent is to cause

files that have experienced I/O activity in the recent past to be relocated to higher performing, perhaps more failure tolerant storage, `ACCAGE` is too coarse a filter. For example, in a rule specifying that files on `tier2` volumes that have been accessed within the last three days should be relocated to `tier1` volumes, no distinction is made between a file that was browsed by a single user and a file that actually was used intensively by applications.

DST implements the concept of I/O temperature and access temperature to overcome these deficiencies. A file's I/O temperature is equal to the number of bytes transferred to or from it over a specified period of time divided by the size of the file. For example, if a file occupies one megabyte of storage at the time of an `fsppadm enforce` operation and the data in the file has been completely read or written 15 times within the last three days, VxFS calculates its 3-day average I/O temperature to be 5 (15 MB of I/O ÷ 1 MB file size ÷ 3 days).

Similarly, a file's average access temperature is the number of read or write requests made to it over a specified number of 24-hour periods divided by the number of periods. Unlike I/O temperature, access temperature is unrelated to file size. A large file to which 20 I/O requests are made over a 2-day period has the same average access temperature as a small file accessed 20 times over a 2-day period.

If a file system's active placement policy includes any `<IOTEMP>` or `<ACCESSTEMP>` clauses, VxFS begins policy enforcement by using information in the file system's FCL file to calculate average I/O activity against all files in the file system during the longest `<PERIOD>` specified in the policy. Shorter specified periods are ignored. VxFS uses these calculations to qualify files for I/O temperature-based relocation and deletion.

See [“About the File Change Log file”](#) on page 114.

Note: If FCL is turned off, I/O temperature-based relocation will not be accurate. When you invoke the `fsppadm enforce` command, the command displays a warning if the FCL is turned off.

As its name implies, the File Change Log records information about changes made to files in a VxFS file system. In addition to recording creations, deletions, extensions, the FCL periodically captures the cumulative amount of I/O activity (number of bytes read and written) on a file-by-file basis. File I/O activity is recorded in the FCL each time a file is opened or closed, as well as at timed intervals to capture information about files that remain open for long periods.

If a file system's active file placement policy contains `<IOTEMP>` clauses, execution of the `fsppadm enforce` command begins with a scan of the FCL to extract I/O activity information over the period of interest for the policy. The period of interest

is the interval between the time at which the `fsppadm enforce` command was issued and that time minus the largest interval value specified in any `<PERIOD>` element in the active policy.

For files with I/O activity during the largest interval, VxFS computes an approximation of the amount of read, write, and total data transfer (the sum of the two) activity by subtracting the I/O levels in the oldest FCL record that pertains to the file from those in the newest. It then computes each file's I/O temperature by dividing its I/O activity by its size at `Tscan`. Dividing by file size is an implicit acknowledgement that relocating larger files consumes more I/O resources than relocating smaller ones. Using this algorithm requires that larger files must have more activity against them in order to reach a given I/O temperature, and thereby justify the resource cost of relocation.

While this computation is an approximation in several ways, it represents an easy to compute, and more importantly, unbiased estimate of relative recent I/O activity upon which reasonable relocation decisions can be based.

File relocation and deletion decisions can be based on read, write, or total I/O activity.

The following XML snippet illustrates the use of `IOTEMP` in a policy rule to specify relocation of low activity files from `tier1` volumes to `tier2` volumes:

```
<RELOCATE>
  <FROM>
    <SOURCE>
      <CLASS>tier1</CLASS>
    </SOURCE>
  </FROM>
  <TO>
    <DESTINATION>
      <CLASS>tier2</CLASS>
    </DESTINATION>
  </TO>
  <WHEN>
    <IOTEMP Type="nrbytes">
      <MAX Flags="lt">3</MAX>
      <PERIOD Units="days">4</PERIOD>
    </IOTEMP>
  </WHEN>
</RELOCATE>
```

This snippet specifies that files to which the rule applies should be relocated from `tier1` volumes to `tier2` volumes if their I/O temperatures fall below 3 over a period of 4 days. The `Type="nrbytes"` XML attribute specifies that total data

transfer activity, which is the the sum of bytes read and bytes written, should be used in the computation. For example, a 50 megabyte file that experienced less than 150 megabytes of data transfer over the 4-day period immediately preceding the `fsppadm enforce` scan would be a candidate for relocation. VxFS considers files that experience no activity over the period of interest to have an I/O temperature of zero. VxFS relocates qualifying files in the order in which it encounters the files in its scan of the file system directory tree.

Using I/O temperature or access temperature rather than a binary indication of activity, such as the POSIX `atime` or `mtime`, minimizes the chance of not relocating files that were only accessed occasionally during the period of interest. A large file that has had only a few bytes transferred to or from it would have a low I/O temperature, and would therefore be a candidate for relocation to `tier2` volumes, even if the activity was very recent.

But, the greater value of I/O temperature or access temperature as a file relocation criterion lies in upward relocation: detecting increasing levels of I/O activity against files that had previously been relocated to lower tiers in a storage hierarchy due to inactivity or low temperatures, and relocating them to higher tiers in the storage hierarchy.

The following XML snippet illustrates relocating files from `tier2` volumes to `tier1` when the activity level against them increases.

```
<RELOCATE>
  <FROM>
    <SOURCE>
      <CLASS>tier2</CLASS>
    </SOURCE>
  </FROM>
  <TO>
    <DESTINATION>
      <CLASS>tier1</CLASS>
    </DESTINATION>
  </TO>
  <WHEN>
    <IOTEMP Type="nrbytes">
      <MAX Flags="gt">5</MAX>
      <PERIOD Units="days">2</PERIOD>
    </IOTEMP>
  </WHEN>
</RELOCATE>
```

The `<RELOCATE>` statement specifies that files on `tier2` volumes whose I/O temperature as calculated using the number of bytes read is above 5 over a 2-day

period are to be relocated to `tier1` volumes. Bytes written to the file during the period of interest are not part of this calculation.

Using I/O temperature rather than a binary indicator of activity as a criterion for file relocation gives administrators a granular level of control over automated file relocation that can be used to attune policies to application requirements. For example, specifying a large value in the `<PERIOD>` element of an upward relocation statement prevents files from being relocated unless I/O activity against them is sustained. Alternatively, specifying a high temperature and a short period tends to relocate files based on short-term intensity of I/O activity against them.

I/O temperature and access temperature utilize the `sqlite3` database for building a temporary table indexed on an inode. This temporary table is used to filter files based on I/O temperature and access temperature. The temporary table is stored in the database file `.__fsppadm_fcliotemp.db`, which resides in the `lost+found` directory of the mount point.

Multiple criteria in file placement policy rule statements

In certain cases, file placement policy rule statements may contain multiple clauses that affect their behavior. In general, when a rule statement contains multiple clauses of a given type, all clauses must be satisfied in order for the statement to be effective. There are four cases of note in which multiple clauses may be used.

Multiple file selection criteria in SELECT statement clauses

Within a single `SELECT` statement, all the selection criteria clauses of a single type are treated as a selection list. A file need only satisfy a single criterion of a given type to be designated.

In the following example, files in any of the `db/datafiles`, `db/indexes`, and `db/logs` directories, all relative to the file system mount point, would be selected:

```
<SELECT>
  <DIRECTORY Flags="nonrecursive">db/datafiles</DIRECTORY>
  <DIRECTORY Flags="nonrecursive">db/indexes</DIRECTORY>
  <DIRECTORY Flags="nonrecursive">db/logs</DIRECTORY>
</SELECT>
```

This example is in direct contrast to the treatment of selection criteria clauses of different types. When a `SELECT` statement includes multiple types of file selection criteria, a file must satisfy one criterion of each type in order for the rule's action statements to apply.

In the following example, a file must reside in one of db/datafiles, db/indexes, or db/logs and be owned by one of DBA_Manager, MFG_DBA, or HR_DBA to be designated for possible action:

```
<SELECT>
  <DIRECTORY Flags="nonrecursive">db/datafiles</DIRECTORY>
  <DIRECTORY Flags="nonrecursive">db/indexes</DIRECTORY>
  <DIRECTORY Flags="nonrecursive">db/logs</DIRECTORY>
  <USER>DBA_Manager</USER>
  <USER>MFG_DBA</USER>
  <USER>HR_DBA</USER>
</SELECT>
```

If a rule includes multiple `SELECT` statements, a file need only satisfy one of them to be selected for action. This property can be used to specify alternative conditions for file selection.

In the following example, a file need only reside in one of db/datafiles, db/indexes, or db/logs or be owned by one of DBA_Manager, MFG_DBA, or HR_DBA to be designated for possible action:

```
<SELECT>
  <DIRECTORY Flags="nonrecursive">db/datafiles</DIRECTORY>
  <DIRECTORY Flags="nonrecursive">db/indexes</DIRECTORY>
  <DIRECTORY Flags="nonrecursive">db/logs</DIRECTORY>
</SELECT>
<SELECT>
  <USER>DBA_Manager</USER>
  <USER>MFG_DBA</USER>
  <USER>HR_DBA</USER>
</SELECT>
```

Multiple placement classes in <ON> clauses of CREATE statements and in <TO> clauses of RELOCATE statements

Both the <ON> clause of the `CREATE` statement and the <TO> clause of the `RELOCATE` statement can specify priority ordered lists of placement classes using multiple <DESTINATION> XML elements. VxFS uses a volume in the first placement class in a list for the designated purpose of file creation or relocation, if possible. If no volume in the first listed class has sufficient free space or if the file system's volume set does not contain any volumes with that placement class, VxFS uses a volume in the second listed class if possible. If no volume in the second listed class can be used, a volume in the third listed class is used if possible, and so forth.

The following example illustrates of three placement classes specified in the `<ON>` clause of a `CREATE` statement:

```
<CREATE>
  <ON>
    <DESTINATION>
      <CLASS>tier1</CLASS>
    </DESTINATION>
    <DESTINATION>
      <CLASS>tier2</CLASS>
    </DESTINATION>
    <DESTINATION>
      <CLASS>tier3</CLASS>
    </DESTINATION>
  </ON>
</CREATE>
```

In this statement, VxFS would allocate space for newly created files designated by the rule's `SELECT` statement on `tier1` volumes if space was available. If no `tier1` volume had sufficient free space, VxFS would attempt to allocate space on a `tier2` volume. If no `tier2` volume had sufficient free space, VxFS would attempt allocation on a `tier3` volume. If sufficient space could not be allocated on a volume in any of the three specified placement classes, allocation would fail with an `ENOSPC` error, even if the file system's volume set included volumes in other placement classes that did have sufficient space.

The `<TO>` clause in the `RELOCATE` statement behaves similarly. VxFS relocates qualifying files to volumes in the first placement class specified if possible, to volumes in the second specified class if not, and so forth. If none of the destination criteria can be met, such as if all specified classes are fully occupied, qualifying files are not relocated, but no error is signaled in this case.

Multiple placement classes in `<FROM>` clauses of `RELOCATE` and `DELETE` statements

The `<FROM>` clause in `RELOCATE` and `DELETE` statements can include multiple source placement classes. However, unlike the `<ON>` and `<TO>` clauses, no order or priority is implied in `<FROM>` clauses. If a qualifying file resides on a volume in any of the placement classes specified in a `<FROM>` clause, it is relocated or deleted regardless of the position of its placement class in the `<FROM>` clause list of classes.

Multiple conditions in <WHEN> clauses of RELOCATE and DELETE statements

The <WHEN> clause in RELOCATE and DELETE statements may include multiple relocation criteria. Any or all of <ACCAGE>, <MODAGE>, <SIZE>, and <IOTEMP> can be specified. When multiple conditions are specified, all must be satisfied in order for a selected file to qualify for relocation or deletion.

In the following example, a selected file would have to be both inactive, that is, not accessed, for more than 30 days and larger than 100 megabytes to be eligible for relocation or deletion:

```
<WHEN>
  <ACCAGE Units="days">
    <MIN Flags="gt">30</MIN>
  </ACCAGE>
  <SIZE Units="MB">
    <MIN Flags="gt">100</MIN>
  </SIZE>
</WHEN>
```

You cannot write rules to relocate or delete a single designated set of files if the files meet one of two or more relocation or deletion criteria.

File placement policy rule and statement ordering

You can use the Dynamic Storage Tiering graphical user interface (GUI) to create any of four types of file placement policy documents. Alternatively, you can use a text editor or XML editor to create XML policy documents directly. The GUI places policy rule statements in the correct order to achieve the desired behavior. If you use a text editor, it is your responsibility to order policy rules and the statements in them so that the desired behavior results.

The rules that comprise a placement policy may occur in any order, but during both file allocation and `fsppadm enforce` relocation scans, the first rule in which a file is designated by a `SELECT` statement is the only rule against which that file is evaluated. Thus, rules whose purpose is to supersede a generally applicable behavior for a special class of files should precede the general rules in a file placement policy document.

The following XML snippet illustrates faulty rule placement with potentially unintended consequences:

```
<?xml version="1.0"?>
<!DOCTYPE FILE_PLACEMENT_POLICY SYSTEM "placement.dtd">
```



```

<FILE_PLACEMENT_POLICY Version="5.0">
  <RULE Name="GeneralRule">
    <SELECT>
      <PATTERN>*</PATTERN>
    </SELECT>
    <CREATE>
      <ON>
        <DESTINATION>
          <CLASS>tier2</CLASS>
        </DESTINATION>
      </ON>
    </CREATE>
    other_statements
  </RULE>
  <RULE Name="DatabaseRule">
    <SELECT>
      <PATTERN>*.db</PATTERN>
    </SELECT>
    <CREATE>
      <ON>
        <DESTINATION>
          <CLASS>tier1</CLASS>
        </DESTINATION>
      </ON>
    </CREATE>
    other_statements
  </RULE>
</FILE_PLACEMENT_POLICY>

```

The `GeneralRule` rule specifies that all files created in the file system, designated by `<PATTERN>*</PATTERN>`, should be created on `tier2` volumes. The `DatabaseRule` rule specifies that files whose names include an extension of `.db` should be created on `tier1` volumes. The `GeneralRule` rule applies to any file created in the file system, including those with a naming pattern of `*.db`, so the `DatabaseRule` rule will never apply to any file. This fault can be remedied by exchanging the order of the two rules. If the `DatabaseRule` rule occurs first in the policy document, VxFS encounters it first when determining where to new place files whose names follow the pattern `*.db`, and correctly allocates space for them on `tier1` volumes. For files to which the `DatabaseRule` rule does not apply, VxFS continues scanning the policy and allocates space according to the specification in the `CREATE` statement of the `GeneralRule` rule.

A similar consideration applies to statements within a placement policy rule. VxFS processes these statements in order, and stops processing on behalf of a file when it encounters a statement that pertains to the file. This can result in unintended behavior.

The following XML snippet illustrates a `RELOCATE` statement and a `DELETE` statement in a rule that is intended to relocate if the files have not been accessed in 30 days, and delete the files if they have not been accessed in 90 days:

```
<RELOCATE>
  <TO>
    <DESTINATION>
      <CLASS>tier2</CLASS>
    </DESTINATION>
  </TO>
  <WHEN>
    <ACCAGE Units="days">
      <MIN Flags="gt">30</MIN>
    </ACCAGE>
  </WHEN>
</RELOCATE>
<DELETE>
  <WHEN>
    <ACCAGE Units="days">
      <MIN Flags="gt">90</MIN>
    </ACCAGE>
  </WHEN>
</DELETE>
```

As written with the `RELOCATE` statement preceding the `DELETE` statement, files will never be deleted, because the `<WHEN>` clause in the `RELOCATE` statement applies to all selected files that have not been accessed for at least 30 days. This includes those that have not been accessed for 90 days. VxFS ceases to process a file against a placement policy when it identifies a statement that applies to that file, so the `DELETE` statement would never occur. This example illustrates the general point that `RELOCATE` and `DELETE` statements that specify less inclusive criteria should precede statements that specify more inclusive criteria in a file placement policy document. The GUI automatically produce the correct statement order for the policies it creates.

File placement policies and extending files

In a VxFS file system with an active file placement policy, the placement class on whose volume a file resides is part of its metadata, and is attached when it is created and updated when it is relocated. When an application extends a file, VxFS allocates the incremental space on the volume occupied by the file if possible. If not possible, VxFS allocates the space on another volume in the same placement class. For example, if a file is created on a `tier1` volume and later relocated to a `tier2` volume, extensions to the file that occur before the relocation have space allocated on a `tier1` volume, while those occurring after to the relocation have their space allocated on `tier2` volumes. When a file is relocated, all of its allocated space, including the space acquired by extension, is relocated to `tier2` volumes in this case.

Quick Reference

This appendix includes the following topics:

- [Command summary](#)
- [Online manual pages](#)
- [Creating a VxFS file system](#)
- [Converting a file system to VxFS](#)
- [Mounting a file system](#)
- [Unmounting a file system](#)
- [Displaying information on mounted file systems](#)
- [Identifying file system types](#)
- [Resizing a file system](#)
- [Backing up and restoring a file system](#)
- [Using quotas](#)

Command summary

Symbolic links to all VxFS command executables are installed in the `/opt/VRTS/bin` directory. Add this directory to the end of your `PATH` environment variable to access the commands.

[Table A-1](#) describes the VxFS-specific commands.

Table A-1 VxFS commands

Command	Description
df	Reports the number of free disk blocks and inodes for a VxFS file system.
diskusg	Generates VxFS disk accounting data by user ID.
extendfs	Extends the size of a VxFS file system.
fcladm	Administers VxFS File Change Logs.
ff	Lists file names and inode information for a VxFS file system.
fiostat	Administers file I/O statistics
fsadm	Resizes or defragments a VxFS file system.
fsapadm	Administers VxFS allocation policies.
fscat	Cats a VxFS file system.
fscdsadm	Performs online CDS operations.
fscdsconv	Performs offline CDS migration tasks on VxFS file systems.
fscdstask	Performs various CDS operations.
fsck	Checks and repairs a VxFS file system.
fsckpt_restore	Restores file systems from VxFS Storage Checkpoints.
fsckptadm	Administers VxFS Storage Checkpoints.
fsclustadm	Manages cluster-mounted VxFS file systems.
fsdb	Debugs VxFS file systems.
fsmap	Displays VxFS file system extent information.
fsppadm	Administers VxFS placement policies.
fsppmk	Creates placement policies.
fstag	Creates, deletes, or lists file tags.
fstyp	Returns the type of file system on a specified disk partition.
fsvmap	Maps volumes of VxFS file systems to files.
fsvoladm	Administers VxFS volumes.

Table A-1 VxFS commands (*continued*)

Command	Description
glmconfig	Configures Group Lock Managers (GLM).
glmdump	Reports stuck Group Lock Managers (GLM) locks in a cluster file system.
glmstat	Group Lock Managers (GLM) statistics gathering utility.
mkfs	Constructs a VxFS file system.
mount	Mounts a VxFS file system.
ncheck	Generates path names from inode numbers for a VxFS file system.
newfs	Creates a new VxFS file system.
qioadmin	Administers VxFS Quick I/O for Databases cache.
qiomkfile	Creates a VxFS Quick I/O device file. This functionality is available only with the Veritas Quick I/O for Databases feature.
qiostat	Displays statistics for VxFS Quick I/O for Databases. This functionality is available only with the Veritas Quick I/O for Databases feature.
quot	Summarizes ownership on a VxFS file system.
quotacheck	Checks VxFS file system quota consistency.
rvxdump	Incrementally dumps file systems.
rvxrestore	Restores a file system incrementally.
setext	Sets extent attributes on a file in a VxFS file system.
vxdump	Incrementally dumps file systems.
vxenablef	Enables specific VxFS features.
vxfsconvert	Converts an unmounted file system to VxFS or upgrades a VxFS disk layout version.
vxfsstat	Displays file system statistics.
vxlsino	Looks up VxFS reverse path names.
vxrestore	Restores a file system incrementally.
vxtunefs	Tunes a VxFS file system.
vxumount	Unmounts a VxFS file system.

Table A-1 VxFS commands (continued)

Command	Description
<code>vxupgrade</code>	Upgrades the disk layout of a mounted VxFS file system.

Online manual pages

This release includes the following online manual pages as part of the `VRTSvxfs` package. These are installed in the appropriate directories under `/opt/VRTS/man` (add this to your `MANPATH` environment variable), but does not update the `windex` database. To ensure that new VxFS manual pages display correctly, update the `windex` database after installing `VRTSvxfs`.

See the `catman(1M)` manual page.

[Table A-2](#) describes the VxFS-specific section 1 manual pages.

Table A-2 Section 1 manual pages

Section 1	Description
<code>fiostat</code>	Administers file I/O statistics.
<code>getext</code>	Gets extent attributes for a VxFS file system.
<code>qloadadmin</code>	Administers VxFS Quick I/O for Databases cache. This functionality is available only with the Veritas Quick I/O for Databases feature.
<code>qiomkfile</code>	Creates a VxFS Quick I/O device file. This functionality is available only with the Veritas Quick I/O for Databases feature.
<code>qiostat</code>	Displays statistics for VxFS Quick I/O for Databases. This functionality is available only with the Veritas Quick I/O for Databases feature.
<code>setext</code>	Sets extent attributes on a file in a VxFS file system.

[Table A-3](#) describes the VxFS-specific section 1M manual pages.

Table A-3 Section 1M manual pages

Section 1M	Description
<code>cfsccluster</code>	Configures SFCFS clusters. This functionality is available only with the Veritas Cluster File System product.
<code>cfsgadm</code>	Adds or deletes shared disk groups to/from a cluster configuration. This functionality is available only with the Veritas Cluster File System product.

Table A-3 Section 1M manual pages (*continued*)

Section 1M	Description
<code>cfsmntadm</code>	Adds, deletes, modifies, and sets policy on cluster mounted file systems. This functionality is available only with the Veritas Cluster File System product.
<code>cfsmount</code> , <code>cfsumount</code>	Mounts or unmounts a cluster file system. This functionality is available only with the Veritas Cluster File System product.
<code>df_vxfs</code>	Reports the number of free disk blocks and inodes for a VxFS file system.
<code>extendfs_vxfs</code>	Extends the size of a VxFS file system.
<code>fcladm</code>	Administers VxFS File Change Logs.
<code>ff_vxfs</code>	Lists file names and inode information for a VxFS file system.
<code>fsadm_vxfs</code>	Resizes or reorganizes a VxFS file system.
<code>fsapadm</code>	Administers VxFS allocation policies.
<code>fscat_vxfs</code>	Cats a VxFS file system.
<code>fscdsadm</code>	Performs online CDS operations.
<code>fscdsconv</code>	Performs offline CDS migration tasks on VxFS file systems.
<code>fscdstask</code>	Performs various CDS operations.
<code>fsck_vxfs</code>	Checks and repairs a VxFS file system.
<code>fsckptadm</code>	Administers VxFS Storage Checkpoints.
<code>fsckpt_restore</code>	Restores file systems from VxFS Storage Checkpoints.
<code>fsclustadm</code>	
<code>fsdb_vxfs</code>	Debugs VxFS file systems.
<code>fsmap</code>	Displays VxFS file system extent information.
<code>fsppadm</code>	Administers VxFS placement policies.
<code>fsvmap</code>	Maps volumes of VxFS file systems to files.
<code>fsvoladm</code>	Administers VxFS volumes.
<code>glmconfig</code>	Configures Group Lock Managers (GLM). This functionality is available only with the Veritas Cluster File System product.
<code>glmdump</code>	Reports stuck Group Lock Managers (GLM) locks in a cluster file system.

Table A-3 Section 1M manual pages (*continued*)

Section 1M	Description
mkfs_vxfs	Constructs a VxFS file system.
mount_vxfs	Mounts a VxFS file system.
ncheck_vxfs	Generates path names from inode numbers for a VxFS file system.
newfs_vxfs	Creates a new VxFS file system.
quot	Summarizes ownership on a VxFS file system.
quotacheck_vxfs	Checks VxFS file system quota consistency.
vxdiskusg	Generates VxFS disk accounting data by user ID.
vxdump	Incrementally dumps file systems.
vxenablef	Enables specific VxFS features.
vxfsconvert	Converts an unmounted file system to VxFS or upgrades a VxFS disk layout version.
vxfsstat	Displays file system statistics.
vxlsino	Looks up VxFS reverse path names.
vxrestore	Restores a file system incrementally.
vxtunefs	Tunes a VxFS file system.
vxupgrade	Upgrades the disk layout of a mounted VxFS file system.
vxumount	Unmounts a VxFS file system.

[Table A-4](#) describes the VxFS-specific section 3 manual pages.

Table A-4 Section 3 manual pages

Section 3	Description
fsckpt_cntl	Performs control functions and Storage Checkpoint state changes.
fsckpt_create	Creates a Storage Checkpoint associated with a file system handle.
fsckpt_createall	Creates a Storage Checkpoint associated with a list of file system handles.
fsckpt_fbmap	Retrieves the block map from a Storage Checkpoint file.
fsckpt_fclose	Closes a Storage Checkpoint file.

Table A-4 Section 3 manual pages (*continued*)

Section 3	Description
<code>fsckpt_finfo</code>	Returns the status information from a Storage Checkpoint file.
<code>fsckpt_fopen</code>	Opens a Storage Checkpoint file.
<code>fsckpt_fpromote</code>	Promotes a file from a Storage Checkpoint into another fileset.
<code>fsckpt_fsclose</code>	Closes a mount point opened for Storage Checkpoint management.
<code>fsckpt_fsopen</code>	Opens a mount point for Storage Checkpoint management.
<code>fsckpt_info</code>	Returns status information on a Storage Checkpoint.
<code>fsckpt_intro</code>	Introduces the VxFS file system Storage Checkpoint API.
<code>fsckpt_mkprimary</code>	Makes a Storage Checkpoint in a VxFS file system the primary fileset for that file system.
<code>fsckpt_opts_create</code>	Creates a Storage Checkpoint associated with a file system handle.
<code>fsckpt_opts_createall</code>	Creates a Storage Checkpoint associated with a list of file system handles.
<code>fsckpt_remove</code>	Removes a Storage Checkpoint from a file system handle.
<code>fsckpt_rename</code>	Renames a Storage Checkpoint from a file system handle.
<code>fsckpt_sa_info</code>	Returns status information on a Storage Checkpoint within a VxFS file system on a block special device.
<code>fsckpt_setqlimit</code> <code>fsckpt_getqlimit</code>	Sets or gets Storage Checkpoint quota limits.
<code>fsckpt_setquota</code>	Turns Storage Checkpoint quotas on or off.
<code>vxfs_ap_alloc2</code>	Allocates an <code>fsap_info2</code> structure.
<code>vxfs_ap_assign_ckpt</code>	Assigns an allocation policy to file data and metadata in a Storage Checkpoint.
<code>vxfs_ap_assign_ckptchain</code>	Assigns an allocation policy for all of the Storage Checkpoints of a VxFS file system.
<code>vxfs_ap_assign_ckptdef</code>	Assigns a default allocation policy for new Storage Checkpoints of a VxFS file system.
<code>vxfs_ap_assign_file</code>	Assigns an allocation policy for file data and metadata.
<code>vxfs_ap_assign_file_pat</code>	Assigns a pattern-based allocation policy for a directory.

Table A-4 Section 3 manual pages (*continued*)

Section 3	Description
<code>vxfs_ap_assign_fs</code>	Assigns an allocation policy for all file data and metadata within a specified file system.
<code>vxfs_ap_assign_fs_pat</code>	Assigns an pattern-based allocation policy for a file system.
<code>vxfs_ap_define</code>	Defines a new allocation policy.
<code>vxfs_ap_define2</code>	Defines a new allocation policy.
<code>vxfs_ap_enforce_ckpt</code>	Reorganizes blocks in a Storage Checkpoint to match a specified allocation policy.
<code>vxfs_ap_enforce_ckptchain</code>	Enforces the allocation policy for all of the Storage Checkpoints of a VxFS file system.
<code>vxfs_ap_enforce_file</code>	Ensures that all blocks in a specified file match the file allocation policy.
<code>vxfs_ap_enforce_file2</code>	Reallocates blocks in a file to match allocation policies.
<code>vxfs_ap_enumerate</code>	Returns information about all allocation policies.
<code>vxfs_ap_enumerate2</code>	Returns information about all allocation policies.
<code>vxfs_ap_free2</code>	Frees one or more <code>fsap_info2</code> structures.
<code>vxfs_ap_query</code>	Returns information about a specific allocation policy.
<code>vxfs_ap_query2</code>	Returns information about a specific allocation policy.
<code>vxfs_ap_query_ckpt</code>	Returns information about allocation policies for each Storage Checkpoint.
<code>vxfs_ap_query_ckptdef</code>	Retrieves the default allocation policies for new Storage Checkpoints of a VxFS file system
<code>vxfs_ap_query_file</code>	Returns information about allocation policies assigned to a specified file.
<code>vxfs_ap_query_file_pat</code>	Returns information about the pattern-based allocation policy assigned to a directory.
<code>vxfs_ap_query_fs</code>	Retrieves allocation policies assigned to a specified file system.
<code>vxfs_ap_query_fs_pat</code>	Returns information about the pattern-based allocation policy assigned to a file system.
<code>vxfs_ap_remove</code>	Deletes a specified allocation policy.

Table A-4 Section 3 manual pages (*continued*)

Section 3	Description
<code>vxfs_fcl_sync</code>	Sets a synchronization point in the VxFS File Change Log.
<code>vxfs_fiostats_dump</code>	Returns file and file range I/O statistics.
<code>vxfs_fiostats_getconfig</code>	Gets file range I/O statistics configuration values.
<code>vxfs_fiostats_set</code>	Turns on and off file range I/O statistics and resets statistics counters.
<code>vxfs_get_iooffsets</code>	Obtains VxFS inode field offsets.
<code>vxfs_inotopath</code>	Returns path names for a given inode number.
<code>vxfs_inostat</code>	Gets the file statistics based on the inode number.
<code>vxfs_inotofd</code>	Gets the file descriptor based on the inode number.
<code>vxfs_nattr_check</code> <code>vxfs_nattr_fcheck</code>	Checks for the existence of named data streams.
<code>vxfs_nattr_link</code>	Links to a named data stream.
<code>vxfs_nattr_open</code>	Opens a named data stream.
<code>vxfs_nattr_rename</code>	Renames a named data stream.
<code>vxfs_nattr_unlink</code>	Removes a named data stream.
<code>vxfs_nattr_utimes</code>	Sets access and modification times for named data streams.
<code>vxfs_vol_add</code>	Adds a volume to a multi-volume file system.
<code>vxfs_vol_clearflags</code>	Clears specified flags on volumes in a multi-volume file system.
<code>vxfs_vol_deencapsulate</code>	De-encapsulates a volume from a multi-volume file system.
<code>vxfs_vol_encapsulate</code>	Encapsulates a volume within a multi-volume file system.
<code>vxfs_vol_encapsulate_bias</code>	Encapsulates a volume within a multi-volume file system.
<code>vxfs_vol_enumerate</code>	Returns information about the volumes within a multi-volume file system.
<code>vxfs_vol_queryflags</code>	Queries flags on volumes in a multi-volume file system.
<code>vxfs_vol_remove</code>	Removes a volume from a multi-volume file system.
<code>vxfs_vol_resize</code>	Resizes a specific volume within a multi-volume file system.

Table A-4 Section 3 manual pages (*continued*)

Section 3	Description
<code>vxfs_vol_setflags</code>	Sets specified flags on volumes in a multi-volume file system.
<code>vxfs_vol_stat</code>	Returns free space information about a component volume within a multi-volume file system.

[Table A-5](#) describes the VxFS-specific section 4 manual pages.

Table A-5 Section 4 manual pages

Section 4	Description
<code>fs_vxfs</code>	Provides the format of a VxFS file system volume.
<code>inode_vxfs</code>	Provides the format of a VxFS file system inode.
<code>tunefstab</code>	Describes the VxFS file system tuning parameters table.

[Table A-6](#) describes the VxFS-specific section 7 manual pages.

Table A-6 Section 7 manual pages

Section 7	Description
<code>vxfsio</code>	Describes the VxFS file system control functions.

Creating a VxFS file system

The `mkfs` command creates a VxFS file system by writing to a special character device file. The special character device is a location or character device node of a particular storage device. The `mkfs` command builds a file system with a root directory and a `lost+found` directory.

Before running `mkfs`, you must create the target device. Refer to your operating system documentation for more information. If you are using a logical device (such as a VxVM volume), see the VxVM documentation for instructions on device initialization.

See the `mkfs(1M)` and `mkfs_vxfs(1M)` manual pages.

To create a file system

- ◆ Use the `mkfs` command to create a file system:

```
mkfs [-F vxfs] [-m] [generic_options] [-o specific_options] \
special [size]
```

<code>-F vxfs</code>	Specifies the VxFS file system type.
<code>-m</code>	Displays the command line that was used to create the file system. The file system must already exist. This option enables you to determine the parameters used to construct the file system.
<i>generic_options</i>	Options common to most other file system types.
<code>-o specific_options</code>	Options specific to VxFS.
<code>-o N</code>	Displays the geometry of the file system and does not write to the device.
<code>-o largefiles</code>	Allows users to create files larger than two gigabytes. The default option is <code>largefiles</code> .
<i>special</i>	Specifies the special device file location or character device node of a particular storage device.
<i>size</i>	Specifies the number of 1024-byte sectors in the file system. If <i>size</i> is not specified, <code>mkfs</code> determines the size of the special device.

Example of creating a file system

The following example creates a VxFS file system of 12288 sectors in size on a VxVM volume.

To create a VxFS file system

- 1 Create the file system:

```
# mkfs -F vxfs /dev/vx/rdisk/diskgroup/volume 12288
version 7 layout
262144 sectors, 262144 blocks of size 1024, log size 1024 blocks
largefiles supported
```

- 2 Mount the newly created file system.

Converting a file system to VxFS

The `vxfsconvert` command can be used to convert a HFS file system to a VxFS file system.

See the `vxfsconvert(1M)` manual page.

To convert a HFS file system to a VxFS file system

- ◆ Use the `vxfsconvert` command to convert a HFS file system to VxFS:

```
vxfsconvert [-l logsize] [-s size] [-efnNvyY] special
```

<code>-e</code>	Estimates the amount of space required to complete the conversion.
<code>-f</code>	Displays the list of supported file system types.
<code>-l logsize</code>	Specifies the size of the file system intent log.
<code>-n N</code>	Assumes a no response to all questions asked by <code>vxfsconvert</code> .
<code>-s siz</code>	Directs <code>vxfsconvert</code> to use free disk space past the current end of the file system to store VxFS metadata.
<code>-v</code>	Specifies verbose mode.
<code>-y Y</code>	Assumes a yes response to all questions asked by <code>vxfsconvert</code> .
<i>special</i>	Specifies the name of the character (raw) device that contains the file system to convert.

Example of converting a file system

The following example converts a HFS file system to a VxFS file system with an intent log size of 4096 blocks.

To convert a HFS file system to a VxFS file system

- ◆ Convert the file system:

```
# vxfsconvert -l 4096 /dev/vx/rdisk/diskgroup/volume
```

Mounting a file system

You can mount a VxFS file system by using the `mount` command. If you enter this command, the generic `mount` command parses the arguments and the `-F fstype` option executes the `mount` command specific to that file system type. For VxFS

and Veritas-installed products, the `generic mount` command executes the VxFS `mount` command from the directory `/sbin/fs/vxfs`. If the `-F` option is not supplied, the command searches the file `/etc/fstab` for a file system and an `fstype` matching the special file or mount point provided. If no file system type is specified, `mount` uses the default file system type (VxFS).

To mount a file system

- ◆ Use the `mount` command to mount a file system:

```
mount [-F vxfs] [generic_options] [-r] [-o specific_options] \  
special mount_point
```

<code>vxfs</code>	File system type.
<code>generic_options</code>	Options common to most other file system types.
<code>specific_options</code>	Options specific to VxFS.
<code>-o ckpt=ckpt_name</code>	Mounts a Storage Checkpoint.
<code>-o cluster</code>	Mounts a file system in shared mode. Available only with the VxFS cluster file system feature.
<code>special</code>	A VxFS block special device.
<code>mount_point</code>	Directory on which to mount the file system.
<code>-r</code>	Mounts the file system as read-only.

Mount options

The `mount` command has numerous options to tailor a file system for various functions and environments.

The following table lists some of the *specific_options*:

Security feature	If security is important, use <code>blkclear</code> to ensure that deleted files are completely erased before the space is reused.
Support for large files	If you specify the <code>largefiles</code> option, you can create files larger than two gigabytes on the file system. The default option is <code>largefiles</code> .

Support for cluster file systems	If you specify the <code>cluster</code> option, the file system is mounted in shared mode. Cluster file systems depend on several other Veritas products that must be correctly configured before a complete clustering environment is enabled.
Using Storage Checkpoints	The <code>ckpt=checkpoint_name</code> option mounts a Storage Checkpoint of a mounted file system that was previously created by the <code>fsckptadm</code> command.
Using databases	If you are using databases with VxFS and if you have installed a license key for the Veritas Quick I/O for Databases feature, the <code>mount</code> command enables Quick I/O by default (the same as specifying the <code>qio</code> option). The <code>noqio</code> option disables Quick I/O. If you do not have Quick I/O, <code>mount</code> ignores the <code>qio</code> option. Alternatively, you can increase database performance using the <code>mount</code> option <code>convosync=direct</code> , which utilizes direct I/O.
News file systems	If you are using <code>cnews</code> , use <code>delaylog</code> (or <code>tmplog</code>), <code>mincache=closesync</code> because <code>cnews</code> does an <code>fsync()</code> on each news file before marking it received. The <code>fsync()</code> is performed synchronously as required, but other options are delayed.
Temporary file systems	For a temporary file system such as <code>/tmp</code> , where performance is more important than data integrity, use <code>tmplog,mincache=tmpcache</code> .
Locking a file system	If you specify the <code>mntlock</code> option, you can lock a file system to disallow unmounting the file system except if the <code>mntunlock</code> option is specified. The <code>mntlock</code> is useful for applications for which you do not want the file systems that the applications are monitoring to be improperly unmounted by other applications or administrators.

See [“Mounting a VxFS file system”](#) on page 34.

See the `fstab(4)`, `fsckptadm(1M)`, `mount(1M)`, and `mount_vxfs(1M)` manual pages.

Example of mounting a file system

The following example mounts the file system `/dev/vx/dsk/fsvol/voll` on the `/ext` directory with read/write access and delayed logging.

To mount the file system

- ◆ Mount the file system:

```
# mount -F vxfs -o delaylog /dev/vx/dsk/fsvol/vol1 /ext
```

Editing the fstab file

You can edit the `/etc/fstab` file to mount a file system automatically at boot time.

You must specify the following:

- The special block device name to mount
- The mount point
- The file system type (vxfs)
- The mount options
- The backup frequency
- Which `fsck` pass looks at the file system

Each entry must be on a single line.

See the `fstab(4)` manual page.

The following is a typical `fstab` file with the new file system on the last line:

```
# System /etc/fstab file.  Static
# information about the file systems
# See fstab(4) and sam(1M) for further
# details on configuring devices.
/dev/vg00/lvol3      /          vxfs delaylog 0 1
/dev/vg00/lvol1      /stand hfs  defaults 0 1
/dev/vg00/lvol4      /tmp       vxfs delaylog 0 2
/dev/vg00/lvol5      /home      vxfs delaylog 0 2
/dev/vg00/lvol6      /opt       vxfs delaylog 0 2
/dev/vg00/lvol7      /usr       vxfs delaylog 0 2
/dev/vg00/lvol8      /var       vxfs delaylog 0 2
/dev/vx/dsk/fsvol    /ext       vxfs delaylog 0 2
```

Unmounting a file system

Use the `umount` command to unmount a currently mounted file system.

See the `vxumount(1M)` manual page.

To unmount a file system

- ◆ Use the `umount` command to unmount a file system:

```
vxumount [-o [force]] mount_point
vxumount [-f] mount_point

vxumount [-o [force]] {special|mount_point}
```

Specify the file system to be unmounted as a *mount_point* or *special*. *special* is the VxFS block special device on which the file system resides.

Example of unmounting a file system

The following are examples of unmounting file systems.

To unmount the file system `/dev/vx/dsk/fsvol/vol1`

- ◆ Unmount the file system:

```
# umount /dev/vx/dsk/fsvol/vol1
```

To unmount all file systems not required by the system

- ◆ Unmount the file system mounted at `/mnt1`:

```
# vxumount /mnt1
```

Displaying information on mounted file systems

Use the `mount` command to display a list of currently mounted file systems.

See the `mount(1M)` and `mount_vxfs(1M)` manual pages.

To view the status of mounted file systems

- ◆ Use the `mount` command to view the status of mounted file systems:

```
mount -v
```

This shows the file system type and `mount` options for all mounted file systems. The `-v` option specifies verbose mode.

Example of displaying information on mounted file systems

The following example shows the result of invoking the `mount` command without options.

To display information on mounted file systems

- ◆ Invoke the `mount` command without options:

```
# mount
/dev/vg00/lvol3 on / type vxfs ioerror=mwdisable,delaylog \
Wed Jun 5 3:23:40 2004
/dev/vg00/lvol8 on /var type vxfs ioerror=mwdisable,delaylog \
Wed Jun 5 3:23:56 2004
/dev/vg00/lvol7 on /usr type vxfs ioerror=mwdisable,delaylog \
Wed Jun 5 3:23:56 2004
/dev/vg00/lvol6 on /tmp type vxfs ioerror=mwdisable,delaylog \
Wed Jun 5 3:23:56 2004
/dev/vg00/lvol5 on /opt type vxfs ioerror=mwdisable,delaylog \
Wed Jun 5 3:23:57 2004
/dev/vg00/lvol11 on /stand type hfs defaults on \
Thu Jun 6 4:17:20 2004
/dev/vgdb/lvol13 on /oracle type vxfs \
ioerror=mwdisable,delaylog Thu Jun 6 4:17:20 2004
/dev/vg00/lvol14 on /home type vxfs \
ioerror=mwdisable,delaylog on Thu Jun 6 4:17:20 2004
/dev/vgdb/lvol19 on /bench type vxfs \
ioerror=mwdisable,delaylog on Thu Jun 6 4:17:11 2004
```

Identifying file system types

Use the `fstyp` command to determine the file system type for a specified file system. This is useful when a file system was created elsewhere and you want to know its type.

See the `fstyp(1M)` manual page.

To determine a file system's type

- ◆ Use the `fstyp` command to determine a file system's type:

```
fstyp -v special
```

special The character (raw) device.

`-v` Specifies verbose mode.

Example of determining a file system's type

The following example uses the `fstyp` command to determine the file system type of the `/dev/vx/dsk/fsvol/vol1` device.

To determine the file system's type

- ◆ Use the `fstyp` command to determine the file system type of the device `/dev/vx/dsk/fsvol/vol1`:

```
# fstyp -v /dev/vx/dsk/fsvol/vol1
```

The output indicates that the file system type is `vxfs`, and displays file system information similar to the following:

```
vxfs  version: 6  f_bsize: 8192  f_fsize: 1024  f_blocks: 1027432  f_bfr
```

Resizing a file system

You can extend or shrink mounted VxFS file systems using the `fsadm` command. Use the `extendfs` command to extend the size of an unmounted file system. A file system using the Version 4 disk layout can be up to two terabytes in size. A file system using the Version 5 disk layout can be up to 32 terabytes in size. A file system using the Version 6 or 7 disk layout can be up to 8 exabytes in size. The size to which a Version 5, 6, or 7 disk layout file system can be increased depends on the file system block size.

See [“About disk layouts”](#) on page 255.

See the `extendfs(1M)` and `fsadm_vxfs(1M)` manual pages.

Extending a file system using `fsadm`

If a VxFS file system is not large enough, you can increase its size. The size of the file system is specified in units of 1024-byte blocks (or sectors).

Note: If a file system is full, busy, or too fragmented, the resize operation may fail.

The device must have enough space to contain the larger file system.

See the `format(1M)` manual page.

See the *Veritas Volume Manager Administrator's Guide*.

To extend a VxFS file system

- ◆ Use the `fsadm` command to extend a VxFS file system:

```
/usr/lib/fs/vxfs/fsadm [-F vxfs] [-b newsize] [-r rawdev] \
mount_point
```

<code>vxfs</code>	The file system type.
<code>newsize</code>	The size (in sectors) to which the file system will increase.
<code>mount_point</code>	The file system's mount point.
<code>-r rawdev</code>	Specifies the path name of the raw device if there is no entry in <code>/etc/fstab</code> and <code>fsadm</code> cannot determine the raw device.

Example of extending a file system

The following is an example of extending a file system with the `fsadm` command.

To extend a file system

- ◆ Extend the VxFS file system mounted on `/ext` to 22528 sectors:

```
# fsadm -F vxfs -b 22528 /ext
```

Shrinking a file system

You can decrease the size of the file system using `fsadm`, even while the file system is mounted.

Note: If a file system is full, busy, or too fragmented, the resize operation may fail.

To decrease the size of a VxFS file system

- ◆ Use the `fsadm` command to decrease the size of a VxFS file system:

```
fsadm [-F vxfs] [-b newsize] [-r rawdev] mount_point
```

<code>vxfs</code>	The file system type.
<code>newsize</code>	The size (in sectors) to which the file system will shrink.
<code>mount_point</code>	The file system's mount point.

`-r rawdev`

Specifies the path name of the raw device if there is no entry in `/etc/fstab` and `fsadm` cannot determine the raw device.

Example of shrinking a file system

The following example shrinks a VxFS file system mounted at `/ext` to 20480 sectors.

To shrink a VxFS file system

- ◆ Shrink a VxFS file system mounted at `/ext` to 20480 sectors:

```
# fsadm -F vxfs -b 20480 /ext
```

Warning: After this operation, there is unused space at the end of the device. You can then resize the device, but be careful not to make the device smaller than the new size of the file system.

Reorganizing a file system

You can reorganize or compact a fragmented file system using `fsadm`, even while the file system is mounted. This may help shrink a file system that could not previously be decreased.

Note: If a file system is full or busy, the reorg operation may fail.

To reorganize a VxFS file system

- ◆ Use the `fsadm` command to reorganize a VxFS file system:

```
fsadm [-F vxfs] [-e] [-d] [-E] [-D] [-r rawdev] mount_point
```

<code>vxfs</code>	The file system type.
<code>-d</code>	Reorders directory entries to put subdirectory entries first, then all other entries in decreasing order of time of last access. Also compacts directories to remove free space.
<code>-D</code>	Reports on directory fragmentation.
<code>-e</code>	Minimizes file system fragmentation. Files are reorganized to have the minimum number of extents.

<code>-E</code>	Reports on extent fragmentation.
<code>mount_point</code>	The file system's mount point.
<code>-r rawdev</code>	Specifies the path name of the raw device if there is no entry in <code>/etc/fstab</code> and <code>fsadm</code> cannot determine the raw device.

Example of reorganizing a file system

The following example reorganizes the file system mounted at `/ext`.

To reorganize a VxFS file system

- ◆ Reorganize the VxFS file system mounted at `/ext`:

```
# fsadm -F vxfs -EeDd /ext
```

Extending a file system using `extendfs`

If a VxFS file system is not mounted, you can use the `extendfs` command to increase the size of the file system.

To extend a VxFS file system

- ◆ Use the `extendfs` command to extend a VxFS file system:

```
extendfs [-F vxfs] [-q] [-v] [-s size] special
```

<code>vxfs</code>	The file system type
<code>-q</code>	Displays the size of special without resizing it
<code>-v</code>	Specifies verbose mode
<code>-s size</code>	Specifies the number of blocks to add to the file system (maximum if not specified)
<code>special</code>	Either a logical volume or a disk partition

Note: If the file system resides on a volume set, the `extendfs` command fails. Use the `fsvoladm` command to extend a multi-volume file system. See the `fsvoladm(1M)` manual page.

Example of extending a VxFS file system

The following example extends a VxFS file system on a VxVM volume.

To increase the capacity of a file system

- 1 Unmount the file system:

```
# umount /dev/vg00/lvol17
```

- 2 Extend the volume so that the volume can contain the larger file system:

```
# lvextend -L larger_size /dev/vg00/lvol17
```

- 3 Extend the file system:

```
# extendfs -F vxfs /dev/vg00/r1vol17
```

- 4 Mount the file system:

```
# mount -F vxfs /dev/vg00/lvol17 mount_point
```

Backing up and restoring a file system

To back up a VxFS file system, you first create a read-only snapshot file system, then back up the snapshot. This procedure lets you keep the main file system on line. The snapshot is a copy of the snapped file system that is frozen at the moment the snapshot is created.

See [“About snapshot file systems”](#) on page 75.

See the `mount(1M)`, `mount_vxfs(1M)`, `vxdump(1M)`, and `vxrestore(1M)` manual pages.

Creating and mounting a snapshot file system

The first step in backing up a VxFS file system is to create and mount a snapshot file system.

To create and mount a snapshot of a VxFS file system

- ◆ Use the `mount` command to create and mount a snapshot of a VxFS file system:

```
mount [-F vxfs] -o snapof=source,[snapsize=size] \  
destination snap_mount_point
```

<i>source</i>	The special device name or mount point of the file system to copy.
<i>destination</i>	The name of the special device on which to create the snapshot.
<i>size</i>	The size of the snapshot file system in sectors.
<i>snap_mount_point</i>	Location where to mount the snapshot; <i>snap_mount_point</i> must exist before you enter this command.

Example of creating and mounting a snapshot of a VxFS file system

The following example creates a snapshot file system of the file system at `/home` on `/dev/vx/dsk/fsvol/vol1`, and mounts it at `/snapmount`.

To create and mount a snapshot file system of a file system

- ◆ Create a snapshot file system of the file system at `/home` on `/dev/vx/dsk/fsvol/vol1` and mount it at `/snapmount`:

```
# mount -F vxfs -o snapof=/home, \
snapsize=32768 /dev/vx/dsk/fsvol/vol1 /snapmount
```

You can now back up the file system.

Backing up a file system

After creating a snapshot file system, you can use `vxdump` to back it up.

To back up a VxFS snapshot file system

- ◆ Use the `vxdump` command to back up a VxFS snapshot file system:

```
vxdump [-c] [-f backupdev] snap_mount_point
```

<code>-c</code>	Specifies using a cartridge tape device.
<i>backupdev</i>	The device on which to back up the file system.
<i>snap_mount_point</i>	The snapshot file system's mount point.

Example of backing up a file system

The following example backs up the VxFS snapshot file system mounted at `/snapmount` to the tape drive with device name `/dev/rmt`.

To back up a VxFS snapshot file system

- ◆ Back up the VxFS snapshot file system mounted at `/snapmount` to the tape drive with device name `/dev/rmt/`:

```
# vxdump -cf /dev/rmt /snapmount
```

Restoring a file system

After backing up the file system, you can restore it using the `vxrestore` command. First, create and mount an empty file system.

To restore a VxFS snapshot file system

- ◆ Use the `vxrestore` command to restore a VxFS snapshot file system:

```
vxrestore [-v] [-x] [filename]
```

`-v` Specifies verbose mode.

`-x` Extracts the named files from the tape.

filename The file or directory to restore. If *filename* is omitted, the root directory, and thus the entire tape, is extracted.

Example of restoring a file system

The following example restores a VxFS snapshot file system from the tape:

To restore a VxFS snapshot file system

- ◆ Restore a VxFS snapshot file system from the tape `/dev/st1` into the mount point `/restore`:

```
# cd /restore
# vxrestore -v -x -f /dev/st1
```

Using quotas

You can use quotas to allocate per-user quotas on VxFS file systems.

See [“Using quotas”](#) on page 110.

See the `edquota(1M)`, `quota(1M)`, `quotaon(1M)`, and `quotaoff(1M)` manual pages.

Turning on quotas

You can enable quotas at mount time or after a file system is mounted. The root directory of the file system must contain a file named `quotas` that is owned by `root`.

To turn on quotas

- 1 Turn on quotas for a mounted file system:

```
quotaon mount_point
```

- 2 Mount a file system and turn on quotas at the same time:

```
mount -F vxfs -o quota special mount_point
```

If the root directory does not contain a `quotas` file, the `mount` command succeeds, but quotas are not turned on.

Example of turning on quotas for a mounted file system

The following example creates a `quotas` file and turns on quotas for a VxFS file system mounted at `/mnt`.

To turn on quotas for a mounted file system

- ◆ Create a `quotas` file if it does not already exist and turn on quotas for a VxFS file system mounted at `/mnt`:

```
# touch /mnt/quotas
# quotaon /mnt
```

Example of turning on quotas at mount time

The following example turns on quotas when the `/dev/vx/dsk/fsvol/voll` file system is mounted.

To turn on quotas for a file system at mount time

- ◆ Turn on quotas at mount time by specifying the `-o quota` option:

```
# mount -F vxfs -o quota /dev/vx/dsk/fsvol/voll /mnt
```

Setting up user quotas

You can set user quotas with the `edquota` command if you have superuser privileges. User quotas can have a soft limit and hard limit. You can modify the

limits or assign them specific values. Users are allowed to exceed the soft limit, but only for a specified time. Disk usage can never exceed the hard limit. The default time limit for exceeding the soft limit is seven days on VxFS file systems.

`edquota` creates a temporary file for a specified user. This file contains on-disk quotas for each mounted VxFS file system that has a quotas file. The temporary file has one or more lines similar to the following:

```
fs /mnt blocks (soft = 0, hard = 0) inodes (soft=0, hard=0)
fs /mnt1 blocks (soft = 100, hard = 200) inodes (soft=10, hard=20)
```

Quotas do not need to be turned on for `edquota` to work. However, the quota limits apply only after quotas are turned on for a given file system.

`edquota` has an option to modify time limits. Modified time limits apply to the entire file system; you cannot set time limits for an individual user.

To set up user quotas

- 1 Invoke the quota editor:

```
edquota username
```

- 2 Modify the time limit:

```
edquota -t
```

Viewing quotas

The superuser or individual user can view disk quotas and usage on VxFS file systems using the `quota` command. This command displays the user's quotas and disk usage on all mounted VxFS file systems where the quotas file exists. You will see all established quotas regardless of whether or not the quotas are actually turned on.

To view quotas for a specific user

- ◆ Use the quota command to view quotas for a specific user:

```
quota -v username
```

Turning off quotas

You can turn off quotas for a mounted file system using the `quotaoff` command.

To turn off quotas for a file system

- ◆ Turn off quotas for a file system:

```
quotaoff mount_point
```


Diagnostic messages

This appendix includes the following topics:

- [File system response to problems](#)
- [About kernel messages](#)
- [Kernel messages](#)
- [About unique message identifiers](#)
- [Unique message identifiers](#)

File system response to problems

When the file system encounters problems, it responds in one of the following ways:

- | | |
|------------------------|--|
| Marking an inode bad | Inodes can be marked bad if an inode update or a directory-block update fails. In these types of failures, the file system does not know what information is on the disk, and considers all the information that it finds to be invalid. After an inode is marked bad, the kernel still permits access to the file name, but any attempt to access the data in the file or change the inode fails. |
| Disabling transactions | If the file system detects an error while writing the intent log, it disables transactions. After transactions are disabled, the files in the file system can still be read or written, but no block or inode frees or allocations, structural changes, directory entry changes, or other changes to metadata are allowed. |

Disabling a file system If an error occurs that compromises the integrity of the file system, VxFS disables itself. If the intent log fails or an inode-list error occurs, the super-block is ordinarily updated (setting the `VX_FULLFSCK` flag) so that the next `fsck` does a full structural check. If this super-block update fails, any further changes to the file system can cause inconsistencies that are undetectable by the intent log replay. To avoid this situation, the file system disables itself.

Recovering a disabled file system

When the file system is disabled, no data can be written to the disk. Although some minor file system operations still work, most simply return `EIO`. The only thing that can be done when the file system is disabled is to do a `umount` and run a full `fsck`.

Although a log replay may produce a clean file system, do a full structural check to be safe.

The file system usually becomes disabled because of disk errors. Disk failures that disable a file system should be fixed as quickly as possible.

See the `fsck_vxfs(1M)` manual page.

To execute a full structural check

- ◆ Use the `fsck` command to execute a full structural check:

```
# fsck -F vxfs -o full -y /dev/vx/rdisk/diskgroup/volume
```

Warning: Be careful when running this command. By specifying the `-y` option, all `fsck` user prompts are answered with a “yes”, which can make irreversible changes if it performs a full file system check.

About kernel messages

Kernel messages are diagnostic or error messages generated by the Veritas File System (VxFS) kernel. Each message has a description and a suggestion on how to handle or correct the underlying problem.

About global message IDs

When a VxFS kernel message displays on the system console, it is preceded by a numerical ID shown in the `msgcnt` field. This ID number increases with each

instance of the message to guarantee that the sequence of events is known when analyzing file system problems.

Each message is also written to an internal kernel buffer that you can view in the file `/var/adm/syslog/syslog.log`.

In some cases, additional data is written to the kernel buffer. For example, if an inode is marked bad, the contents of the bad inode are written. When an error message is displayed on the console, you can use the unique message ID to find the message in `/var/adm/syslog/syslog.log` and obtain the additional information.

Kernel messages

Some commonly encountered kernel messages are described on the following table:

Table B-1 Kernel messages

Message Number	Message and Definition
001	<p>NOTICE: msgcnt x: mesg 001: V-2-1: vx_nospace - <i>mount_point</i> file system full (n block extent)</p> <ul style="list-style-type: none">■ Description The file system is out of space. Often, there is plenty of space and one runaway process used up all the remaining free space. In other cases, the available free space becomes fragmented and unusable for some files.■ Action Monitor the free space in the file system and prevent it from becoming full. If a runaway process has used up all the space, stop that process, find the files created by the process, and remove them. If the file system is out of space, remove files, defragment, or expand the file system. To remove files, use the <code>find</code> command to locate the files that are to be removed. To get the most space with the least amount of work, remove large files or file trees that are no longer needed. To defragment or expand the file system, use the <code>fsadm</code> command. See the <code>fsadm_vxfs(1M)</code> manual page.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
002	<p>WARNING: msgcnt x: msg 002: V-2-2: vx_snap_strategy - <i>mount_point</i> file system write attempt to read-only file system</p> <p>WARNING: msgcnt x: msg 002: V-2-2: vx_snap_copyblk - <i>mount_point</i> file system write attempt to read-only file system</p> <ul style="list-style-type: none"> ■ Description The kernel tried to write to a read-only file system. This is an unlikely problem, but if it occurs, the file system is disabled. ■ Action The file system was not written, so no action is required. Report this as a bug to your customer support organization.
003, 004, 005	<p>WARNING: msgcnt x: msg 003: V-2-3: vx_mapbad - <i>mount_point</i> file system free extent bitmap in au aun marked bad</p> <p>WARNING: msgcnt x: msg 004: V-2-4: vx_mapbad - <i>mount_point</i> file system free inode bitmap in au aun marked bad</p> <p>WARNING: msgcnt x: msg 005: V-2-5: vx_mapbad - <i>mount_point</i> file system inode extended operation bitmap in au aun marked bad</p> <ul style="list-style-type: none"> ■ Description If there is an I/O failure while writing a bitmap, the map is marked bad. The kernel considers the maps to be invalid, so does not do any more resource allocation from maps. This situation can cause the file system to report out of space or out of inode error messages even though df may report an adequate amount of free space. This error may also occur due to bitmap inconsistencies. If a bitmap fails a consistency check, or blocks are freed that are already free in the bitmap, the file system has been corrupted. This may have occurred because a user or process wrote directly to the device or used <code>fsdb</code> to change the file system. The <code>VX_FULLFSCK</code> flag is set. If the map that failed was a free extent bitmap, and the <code>VX_FULLFSCK</code> flag cannot be set, then the file system is disabled. ■ Action Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process was writing to the device, report the problem to your customer support organization. Unmount the file system and use <code>fsck</code> to run a full structural check.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
006, 007	<p>WARNING: msgcnt <i>x</i>: mesg 006: V-2-6: vx_sumupd - <i>mount_point</i> file system summary update in au aun failed</p> <p>WARNING: msgcnt <i>x</i>: mesg 007: V-2-7: vx_sumupd - <i>mount_point</i> file system summary update in inode au iaun failed</p> <p>■ Description An I/O error occurred while writing the allocation unit or inode allocation unit bitmap summary to disk. This sets the <code>VX_FULLFSCK</code> flag on the file system. If the <code>VX_FULLFSCK</code> flag cannot be set, the file system is disabled.</p> <p>■ Action Check the console log for I/O errors. If the problem was caused by a disk failure, replace the disk before the file system is mounted for write access, and use <code>fsck</code> to run a full structural check.</p>
008, 009	<p>WARNING: msgcnt <i>x</i>: mesg 008: V-2-8: vx_direrr: function - <i>mount_point</i> file system dir inode <i>dir_inumber</i> dev/block device_ID/block dirent inode <i>dirent_inumber</i> error <i>errno</i></p> <p>WARNING: msgcnt <i>x</i>: mesg 009: V-2-9: vx_direrr: function - <i>mount_point</i> file system dir inode <i>dir_inumber</i> dirent inode <i>dirent_inumber</i> immediate directory error <i>errno</i></p> <p>■ Description A directory operation failed in an unexpected manner. The mount point, inode, and block number identify the failing directory. If the inode is an immediate directory, the directory entries are stored in the inode, so no block number is reported. If the error is <code>ENOENT</code> or <code>ENOTDIR</code>, an inconsistency was detected in the directory block. This inconsistency could be a bad free count, a corrupted hash chain, or any similar directory structure error. If the error is <code>EIO</code> or <code>ENXIO</code>, an I/O failure occurred while reading or writing the disk block. The <code>VX_FULLFSCK</code> flag is set in the super-block so that <code>fsck</code> will do a full structural check the next time it is run.</p> <p>■ Action Check the console log for I/O errors. If the problem was caused by a disk failure, replace the disk before the file system is mounted for write access. Unmount the file system and use <code>fsck</code> to run a full structural check.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
010	<p>WARNING: msgcnt <i>x</i>: msg 010: V-2-10: vx_ialloc - <i>mount_point</i> file system inode <i>inumber</i> not free</p> <ul style="list-style-type: none"> ■ Description When the kernel allocates an inode from the free inode bitmap, it checks the mode and link count of the inode. If either is non-zero, the free inode bitmap or the inode list is corrupted. The <code>VX_FULLFSCK</code> flag is set in the super-block so that <code>fsck</code> will do a full structural check the next time it is run. ■ Action Unmount the file system and use <code>fsck</code> to run a full structural check.
011	<p>NOTICE: msgcnt <i>x</i>: msg 011: V-2-11: vx_noinode - <i>mount_point</i> file system out of inodes</p> <ul style="list-style-type: none"> ■ Description The file system is out of inodes. ■ Action Monitor the free inodes in the file system. If the file system is getting full, create more inodes either by removing files or by expanding the file system. See the <code>fsadm_vxfs(1M)</code> online manual page.
012	<p>WARNING: msgcnt <i>x</i>: msg 012: V-2-12: vx_iget - <i>mount_point</i> file system invalid inode number <i>inumber</i></p> <ul style="list-style-type: none"> ■ Description When the kernel tries to read an inode, it checks the inode number against the valid range. If the inode number is out of range, the data structure that referenced the inode number is incorrect and must be fixed. The <code>VX_FULLFSCK</code> flag is set in the super-block so that <code>fsck</code> will do a full structural check the next time it is run. ■ Action Unmount the file system and use <code>fsck</code> to run a full structural check.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
013	<p>WARNING: msgcnt <i>x</i>: msg 013: V-2-13: vx_iposition - <i>mount_point</i> file system inode <i>inumber</i> invalid inode list extent</p> <p>■ Description</p> <p>For a Version 2 and above disk layout, the inode list is dynamically allocated. When the kernel tries to read an inode, it must look up the location of the inode in the inode list file. If the kernel finds a bad extent, the inode cannot be accessed. All of the inode list extents are validated when the file system is mounted, so if the kernel finds a bad extent, the integrity of the inode list is questionable. This is a very serious error.</p> <p>The <code>VX_FULFSCCK</code> flag is set in the super-block and the file system is disabled.</p> <p>■ Action</p> <p>Unmount the file system and use <code>fsck</code> to run a full structural check.</p>
014	<p>WARNING: msgcnt <i>x</i>: msg 014: V-2-14: vx_iget - inode table overflow</p> <p>■ Description</p> <p>All the system in-memory inodes are busy and an attempt was made to use a new inode.</p> <p>■ Action</p> <p>Look at the processes that are running and determine which processes are using inodes. If it appears there are runaway processes, they might be tying up the inodes. If the system load appears normal, increase the <code>vx_ninode</code> parameter in the kernel.</p> <p>Note: The tunable parameter <code>vx_ninode</code> is used to set the value of <code>vxfs_ninode</code>.</p> <p>See “Tuning the VxFS file system” on page 44.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
015	<p>WARNING: msgcnt <i>x</i>: msg 015: V-2-15: vx_ibadinactive - <i>mount_point</i> file system cannot mark inode <i>inumber</i> bad</p> <p>WARNING: msgcnt <i>x</i>: msg 015: V-2-15: vx_ilsterr - <i>mount_point</i> file system cannot mark inode <i>inumber</i> bad</p> <p>■ Description</p> <p>An attempt to mark an inode bad on disk, and the super-block update to set the <code>VX_FULFSCCK</code> flag, failed. This indicates that a catastrophic disk error may have occurred since both an inode list block and the super-block had I/O failures. The file system is disabled to preserve file system integrity.</p> <p>■ Action</p> <p>Unmount the file system and use <code>fscck</code> to run a full structural check. Check the console log for I/O errors. If the disk failed, replace it before remounting the file system.</p>
016	<p>WARNING: msgcnt <i>x</i>: msg 016: V-2-16: vx_ilsterr - <i>mount_point</i> file system error reading inode <i>inumber</i></p> <p>■ Description</p> <p>An I/O error occurred while reading the inode list. The <code>VX_FULFSCCK</code> flag is set.</p> <p>■ Action</p> <p>Check the console log for I/O errors. If the problem was caused by a disk failure, replace the disk before the file system is mounted for write access. Unmount the file system and use <code>fscck</code> to run a full structural check.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
017	

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
	WARNING: msgcnt x: msg 017: V-2-17: vx_attr_getblk - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_attr_iget - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_attr_indadd - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_attr_indtrunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_attr_iremove - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_bmap - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_bmap_indirect_ext4 - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_delbuf_flush - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_dio_iovec - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_dirbread - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_dircreate - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_dirlook - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_doextop_iau - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_doextop_now - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_do_getpage - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_enter_ext4 - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_exttrunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core
	WARNING: msgcnt x: msg 017: V-2-17: vx_get_alloc - <i>mount_point</i>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
	file system inode <i>inumber</i> marked bad in core
017 (continued)	<p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_ilisterr - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_indtrunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_iread - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_remove - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_remove_attr - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_logwrite_flush - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_oltmount_iget - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_overlay_bmap - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_readnomap - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_reorg_trunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_stablestore - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_tranitimes - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_trunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_write_alloc2 - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_write_default - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p> <p>WARNING: msgcnt <i>x</i>: mesg 017: V-2-17: vx_zero_alloc - <i>mount_point</i> file system inode <i>inumber</i> marked bad in core</p>

Table B-1 Kernel messages (continued)

Message Number	Message and Definition
017 (continued)	<div><div>■ Description</div><p>When inode information is no longer dependable, the kernel marks it bad in memory. This is followed by a message to mark it bad on disk as well unless the mount command <code>ioerror</code> option is set to <code>disable</code>, or there is subsequent I/O failure when updating the inode on disk. No further operations can be performed on the inode.</p><p>The most common reason for marking an inode bad is a disk I/O failure. If there is an I/O failure in the inode list, on a directory block, or an indirect address extent, the integrity of the data in the inode, or the data the kernel tried to write to the inode list, is questionable. In these cases, the disk driver prints an error message and one or more inodes are marked bad.</p><p>The kernel also marks an inode bad if it finds a bad extent address, invalid inode fields, or corruption in directory data blocks during a validation check. A validation check failure indicates the file system has been corrupted. This usually occurs because a user or process has written directly to the device or used <code>fsdb</code> to change the file system.</p><p>The <code>VX_FULFSCK</code> flag is set in the super-block so <code>fsck</code> will do a full structural check the next time it is run.</p><div><div>■ Action</div><p>Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process is writing to the device, report the problem to your customer support organization. In either case, unmount the file system. The file system can be remounted without a full <code>fsck</code> unless the <code>VX_FULFSCK</code> flag is set for the file system.</p></div></div>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
019	<p>WARNING: msgcnt <i>x</i>: mesg 019: V-2-19: vx_log_add - <i>mount_point</i> file system log overflow</p> <ul style="list-style-type: none">■ Description Log ID overflow. When the log ID reaches <code>VX_MAXLOGID</code> (approximately one billion by default), a flag is set so the file system resets the log ID at the next opportunity. If the log ID has not been reset, when the log ID reaches <code>VX_DISLOGID</code> (approximately <code>VX_MAXLOGID</code> plus 500 million by default), the file system is disabled. Since a log reset will occur at the next 60 second sync interval, this should never happen.■ Action Unmount the file system and use <code>fsck</code> to run a full structural check.
020	<p>WARNING: msgcnt <i>x</i>: mesg 020: V-2-20: vx_logerr - <i>mount_point</i> file system log error <i>errno</i></p> <ul style="list-style-type: none">■ Description Intent log failed. The kernel will try to set the <code>VX_FULLEFCK</code> and <code>VX_LOGBAD</code> flags in the super-block to prevent running a log replay. If the super-block cannot be updated, the file system is disabled.■ Action Unmount the file system and use <code>fsck</code> to run a full structural check. Check the console log for I/O errors. If the disk failed, replace it before remounting the file system.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
021	<p>WARNING: msgcnt x: msg 021: V-2-21: vx_fs_init - <i>mount_point</i> file system validation failure</p> <p>■ Description</p> <p>When a VxFS file system is mounted, the structure is read from disk. If the file system is marked clean, the structure is correct and the first block of the intent log is cleared.</p> <p>If there is any I/O problem or the structure is inconsistent, the kernel sets the <code>VX_FULLFCK</code> flag and the mount fails.</p> <p>If the error is not related to an I/O failure, this may have occurred because a user or process has written directly to the device or used <i>fsdb</i> to change the file system.</p> <p>■ Action</p> <p>Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process is writing to the device, report the problem to your customer support organization. In either case, unmount the file system and use <code>fsck</code> to run a full structural check.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
022	<p>WARNING: msgcnt x: msg 022: V-2-22: vx_mountroot - root file system remount failed</p> <p>■ Description</p> <p>The remount of the root file system failed. The system will not be usable if the root file system cannot be remounted for read/write access.</p> <p>When a root Veritas File System is first mounted, it is mounted for read-only access. After <code>fsck</code> is run, the file system is remounted for read/write access. The remount fails if <code>fsck</code> completed a resize operation or modified a file that was opened before the <code>fsck</code> was run. It also fails if an I/O error occurred during the remount. Usually, the system halts or reboots automatically.</p> <p>■ Action</p> <p>Reboot the system. The system either remounts the root cleanly or runs a full structural <code>fsck</code> and remounts cleanly. If the remount succeeds, no further action is necessary.</p> <p>Check the console log for I/O errors. If the disk has failed, replace it before the file system is mounted for write access.</p> <p>If the system won't come up and a full structural <code>fsck</code> hasn't been run, reboot the system on a backup root and manually run a full structural <code>fsck</code>. If the problem persists after the full structural <code>fsck</code> and there are no I/O errors, contact your customer support organization.</p>
023	<p>WARNING: msgcnt x: msg 023: V-2-23: vx_unmountroot - root file system is busy and cannot be unmounted cleanly</p> <p>■ Description</p> <p>There were active files in the file system and they caused the unmount to fail.</p> <p>When the system is halted, the root file system is unmounted. This happens occasionally when a process is hung and it cannot be killed before unmounting the root.</p> <p>■ Action</p> <p><code>fsck</code> will run when the system is rebooted. It should clean up the file system. No other action is necessary.</p> <p>If the problem occurs every time the system is halted, determine the cause and contact your customer support organization.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
024	<p>WARNING: msgcnt <i>x</i>: msg 024: V-2-24: vx_cutwait - <i>mount_point</i> file system current usage table update error</p> <ul style="list-style-type: none"> ■ Description Update to the current usage table (CUT) failed. For a Version 2 disk layout, the CUT contains a fileset version number and total number of blocks used by each fileset. The <code>VX_FULFSCK</code> flag is set in the super-block. If the super-block cannot be written, the file system is disabled. ■ Action Unmount the file system and use <code>fsck</code> to run a full structural check.
025	<p>WARNING: msgcnt <i>x</i>: msg 025: V-2-25: vx_wsuper - <i>mount_point</i> file system super-block update failed</p> <ul style="list-style-type: none"> ■ Description An I/O error occurred while writing the super-block during a resize operation. The file system is disabled. ■ Action Unmount the file system and use <code>fsck</code> to run a full structural check. Check the console log for I/O errors. If the problem is a disk failure, replace the disk before the file system is mounted for write access.
026	<p>WARNING: msgcnt <i>x</i>: msg 026: V-2-26: vx_snap_copyblk - <i>mount_point</i> primary file system read error</p> <ul style="list-style-type: none"> ■ Description Snapshot file system error. When the primary file system is written, copies of the original data must be written to the snapshot file system. If a read error occurs on a primary file system during the copy, any snapshot file system that doesn't already have a copy of the data is out of date and must be disabled. ■ Action An error message for the primary file system prints. Resolve the error on the primary file system and rerun any backups or other applications that were using the snapshot that failed when the error occurred.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
027	<p>WARNING: msgcnt x: mesg 027: V-2-27: vx_snap_bpcopy - <i>mount_point</i> snapshot file system write error</p> <p>■ Description</p> <p>A write to the snapshot file system failed.</p> <p>As the primary file system is updated, copies of the original data are read from the primary file system and written to the snapshot file system. If one of these writes fails, the snapshot file system is disabled.</p> <p>■ Action</p> <p>Check the console log for I/O errors. If the disk has failed, replace it. Resolve the error on the disk and rerun any backups or other applications that were using the snapshot that failed when the error occurred.</p>
028	<p>WARNING: msgcnt x: mesg 028: V-2-28: vx_snap_alloc - <i>mount_point</i> snapshot file system out of space</p> <p>■ Description</p> <p>The snapshot file system ran out of space to store changes.</p> <p>During a snapshot backup, as the primary file system is modified, the original data is copied to the snapshot file system. This error can occur if the snapshot file system is left mounted by mistake, if the snapshot file system was given too little disk space, or the primary file system had an unexpected burst of activity. The snapshot file system is disabled.</p> <p>■ Action</p> <p>Make sure the snapshot file system was given the correct amount of space. If it was, determine the activity level on the primary file system. If the primary file system was unusually busy, rerun the backup. If the primary file system is no busier than normal, move the backup to a time when the primary file system is relatively idle or increase the amount of disk space allocated to the snapshot file system.</p> <p>Rerun any backups that failed when the error occurred.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
029, 030	<p>WARNING: msgcnt <i>x</i>: msg 029: V-2-29: vx_snap_getbp - <i>mount_point</i> snapshot file system block map write error</p> <p>WARNING: msgcnt <i>x</i>: msg 030: V-2-30: vx_snap_getbp - <i>mount_point</i> snapshot file system block map read error</p> <p>■ Description During a snapshot backup, each snapshot file system maintains a block map on disk. The block map tells the snapshot file system where data from the primary file system is stored in the snapshot file system. If an I/O operation to the block map fails, the snapshot file system is disabled.</p> <p>■ Action Check the console log for I/O errors. If the disk has failed, replace it. Resolve the error on the disk and rerun any backups that failed when the error occurred.</p>
031	<p>WARNING: msgcnt <i>x</i>: msg 031: V-2-31: vx_disable - <i>mount_point</i> file system disabled</p> <p>■ Description File system disabled, preceded by a message that specifies the reason. This usually indicates a serious disk problem.</p> <p>■ Action Unmount the file system and use <code>fsck</code> to run a full structural check. If the problem is a disk failure, replace the disk before the file system is mounted for write access.</p>
032	<p>WARNING: msgcnt <i>x</i>: msg 032: V-2-32: vx_disable - <i>mount_point</i> snapshot file system disabled</p> <p>■ Description Snapshot file system disabled, preceded by a message that specifies the reason.</p> <p>■ Action Unmount the snapshot file system, correct the problem specified by the message, and rerun any backups that failed due to the error.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
033	<p>WARNING: msgcnt <i>x</i>: msg 033: V-2-33: vx_check_badblock - <i>mount_point</i> file system had an I/O error, setting VX_FULLFSCK</p> <p>■ Description</p> <p>When the disk driver encounters an I/O error, it sets a flag in the super-block structure. If the flag is set, the kernel will set the VX_FULLFSCK flag as a precautionary measure. Since no other error has set the VX_FULLFSCK flag, the failure probably occurred on a data block.</p> <p>■ Action</p> <p>Unmount the file system and use <code>fsck</code> to run a full structural check. Check the console log for I/O errors. If the problem is a disk failure, replace the disk before the file system is mounted for write access.</p>
034	<p>WARNING: msgcnt <i>x</i>: msg 034: V-2-34: vx_resetlog - <i>mount_point</i> file system cannot reset log</p> <p>■ Description</p> <p>The kernel encountered an error while resetting the log ID on the file system. This happens only if the super-block update or log write encountered a device failure. The file system is disabled to preserve its integrity.</p> <p>■ Action</p> <p>Unmount the file system and use <code>fsck</code> to run a full structural check. Check the console log for I/O errors. If the problem is a disk failure, replace the disk before the file system is mounted for write access.</p>
035	<p>WARNING: msgcnt <i>x</i>: msg 035: V-2-35: vx_inactive - <i>mount_point</i> file system inactive of locked inode <i>inumber</i></p> <p>■ Description</p> <p>VOP_INACTIVE was called for an inode while the inode was being used. This should never happen, but if it does, the file system is disabled.</p> <p>■ Action</p> <p>Unmount the file system and use <code>fsck</code> to run a full structural check. Report as a bug to your customer support organization.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
036	<p>WARNING: msgcnt x: msg 036: V-2-36: vx_lctbad - <i>mount_point</i> file system link count table <i>lctnumber</i> bad</p> <ul style="list-style-type: none"> ■ Description Update to the link count table (LCT) failed. For a Version 2 and above disk layout, the LCT contains the link count for all the structural inodes. The <code>VX_FULFSCCK</code> flag is set in the super-block. If the super-block cannot be written, the file system is disabled. ■ Action Unmount the file system and use <code>fsck</code> to run a full structural check.
037	<p>WARNING: msgcnt x: msg 037: V-2-37: vx_metaioerr - function - <i>volume_name</i> file system meta data [read write] error in dev/block device_ID/block</p> <ul style="list-style-type: none"> ■ Description A read or a write error occurred while accessing file system metadata. The full <code>fsck</code> flag on the file system was set. The message specifies whether the disk I/O that failed was a read or a write. File system metadata includes inodes, directory blocks, and the file system log. If the error was a write error, it is likely that some data was lost. This message should be accompanied by another file system message describing the particular file system metadata affected, as well as a message from the disk driver containing information about the disk I/O error. ■ Action Resolve the condition causing the disk error. If the error was the result of a temporary condition (such as accidentally turning off a disk or a loose cable), correct the condition. Check for loose cables, etc. Unmount the file system and use <code>fsck</code> to run a full structural check (possibly with loss of data). In case of an actual disk error, if it was a read error and the disk driver remaps bad sectors on write, it may be fixed when <code>fsck</code> is run since <code>fsck</code> is likely to rewrite the sector with the read error. In other cases, you replace or reformat the disk drive and restore the file system from backups. Consult the documentation specific to your system for information on how to recover from disk errors. The disk driver should have printed a message that may provide more information.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
038	<p>WARNING: msgcnt x: msg 038: V-2-38: vx_dataioerr - <i>volume_name</i> file system file data [read write] error in dev/block device_ID/block</p> <p>■ Description</p> <p>A read or a write error occurred while accessing file data. The message specifies whether the disk I/O that failed was a read or a write. File data includes data currently in files and free blocks. If the message is printed because of a read or write error to a file, another message that includes the inode number of the file will print. The message may be printed as the result of a read or write error to a free block, since some operations allocate an extent and immediately perform I/O to it. If the I/O fails, the extent is freed and the operation fails. The message is accompanied by a message from the disk driver regarding the disk I/O error.</p> <p>■ Action</p> <p>Resolve the condition causing the disk error. If the error was the result of a temporary condition (such as accidentally turning off a disk or a loose cable), correct the condition. Check for loose cables, etc. If any file data was lost, restore the files from backups. Determine the file names from the inode number. See the ncheck(1M) manual page.</p> <p>If an actual disk error occurred, make a backup of the file system, replace or reformat the disk drive, and restore the file system from the backup. Consult the documentation specific to your system for information on how to recover from disk errors. The disk driver should have printed a message that may provide more information.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
039	<p>WARNING: msgcnt x: mesg 039: V-2-39: vx_writesuper - file system super-block write error</p> <p>■ Description</p> <p>An attempt to write the file system super block failed due to a disk I/O error. If the file system was being mounted at the time, the mount will fail. If the file system was mounted at the time and the full <code>fsck</code> flag was being set, the file system will probably be disabled and Message 031 will also be printed. If the super-block was being written as a result of a sync operation, no other action is taken.</p> <p>■ Action</p> <p>Resolve the condition causing the disk error. If the error was the result of a temporary condition (such as accidentally turning off a disk or a loose cable), correct the condition. Check for loose cables, etc. Unmount the file system and use <code>fsck</code> to run a full structural check.</p> <p>If an actual disk error occurred, make a backup of the file system, replace or reformat the disk drive, and restore the file system from backups. Consult the documentation specific to your system for information on how to recover from disk errors. The disk driver should have printed a message that may provide more information.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
056	<p>WARNING: msgcnt <i>x</i>: msg 056: V-2-56: vx_mapbad - <i>mount_point</i> file system extent allocation unit state bitmap number <i>number</i> marked bad</p> <p>■ Description</p> <p>If there is an I/O failure while writing a bitmap, the map is marked bad. The kernel considers the maps to be invalid, so does not do any more resource allocation from maps. This situation can cause the file system to report “out of space” or “out of inode” error messages even though <i>df</i> may report an adequate amount of free space.</p> <p>This error may also occur due to bitmap inconsistencies. If a bitmap fails a consistency check, or blocks are freed that are already free in the bitmap, the file system has been corrupted. This may have occurred because a user or process wrote directly to the device or used <i>fsdb</i> to change the file system.</p> <p>The <code>VX_FULFSCCK</code> flag is set. If the <code>VX_FULFSCCK</code> flag cannot be set, the file system is disabled.</p> <p>■ Action</p> <p>Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process was writing to the device, report the problem to your customer support organization. Unmount the file system and use <i>fsck</i> to run a full structural check.</p>
057	<p>WARNING: msgcnt <i>x</i>: msg 057: V-2-57: vx_esum_bad - <i>mount_point</i> file system extent allocation unit summary number <i>number</i> marked bad</p> <p>■ Description</p> <p>An I/O error occurred reading or writing an extent allocation unit summary.</p> <p>The <code>VX_FULFSCCK</code> flag is set. If the <code>VX_FULFSCCK</code> flag cannot be set, the file system is disabled.</p> <p>■ Action</p> <p>Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process was writing to the device, report the problem to your customer support organization. Unmount the file system and use <i>fsck</i> to run a full structural check.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
058	<p>WARNING: msgcnt x: msg 058: V-2-58: vx_isum_bad - <i>mount_point</i> file system inode allocation unit summary number <i>number</i> marked bad</p> <p>■ Description</p> <p>An I/O error occurred reading or writing an inode allocation unit summary.</p> <p>The <code>VX_FULFSCCK</code> flag is set. If the <code>VX_FULFSCCK</code> flag cannot be set, the file system is disabled.</p> <p>■ Action</p> <p>Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process was writing to the device, report the problem to your customer support organization. Unmount the file system and use <code>fsck</code> to run a full structural check.</p>
059	<p>WARNING: msgcnt x: msg 059: V-2-59: vx_snap_getbitbp - <i>mount_point</i> snapshot file system bitmap write error</p> <p>■ Description</p> <p>An I/O error occurred while writing to the snapshot file system bitmap. There is no problem with the snapped file system, but the snapshot file system is disabled.</p> <p>■ Action</p> <p>Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process was writing to the device, report the problem to your customer support organization. Restart the snapshot on an error free disk partition. Rerun any backups that failed when the error occurred.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
060	<p>WARNING: msgcnt <i>x</i>: mesg 060: V-2-60: vx_snap_getbitbp - <i>mount_point</i> snapshot file system bitmap read error</p> <ul style="list-style-type: none">■ Description An I/O error occurred while reading the snapshot file system bitmap. There is no problem with snapped file system, but the snapshot file system is disabled.■ Action Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process was writing to the device, report the problem to your customer support organization. Restart the snapshot on an error free disk partition. Rerun any backups that failed when the error occurred.
061	<p>WARNING: msgcnt <i>x</i>: mesg 061: V-2-61: vx_resize - <i>mount_point</i> file system remount failed</p> <ul style="list-style-type: none">■ Description During a file system resize, the remount to the new size failed. The <code>VX_FULLFSCK</code> flag is set and the file system is disabled.■ Action Unmount the file system and use <code>fsck</code> to run a full structural check. After the check, the file system shows the new size.
062	<p>NOTICE: msgcnt <i>x</i>: mesg 062: V-2-62: vx_attr_creatop - invalid disposition returned by attribute driver</p> <ul style="list-style-type: none">■ Description A registered extended attribute intervention routine returned an invalid return code to the VxFS driver during extended attribute inheritance.■ Action Determine which vendor supplied the registered extended attribute intervention routine and contact their customer support organization.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
063	<p>WARNING: msgcnt <i>x</i>: msg 063: V-2-63: vx_fset_markbad - <i>mount_point</i> file system <i>mount_point</i> fileset (index <i>number</i>) marked bad</p> <p>■ Description An error occurred while reading or writing a fileset structure. <code>VX_FULLFSCK</code> flag is set. If the <code>VX_FULLFSCK</code> flag cannot be set, the file system is disabled.</p> <p>■ Action Unmount the file system and use <code>fsck</code> to run a full structural check.</p>
064	<p>WARNING: msgcnt <i>x</i>: msg 064: V-2-64: vx_ivalidate - <i>mount_point</i> file system inode number version number exceeds fileset's</p> <p>■ Description During inode validation, a discrepancy was found between the inode version number and the fileset version number. The inode may be marked bad, or the fileset version number may be changed, depending on the ratio of the mismatched version numbers. <code>VX_FULLFSCK</code> flag is set. If the <code>VX_FULLFSCK</code> flag cannot be set, the file system is disabled.</p> <p>■ Action Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process is writing to the device, report the problem to your customer support organization. In either case, unmount the file system and use <code>fsck</code> to run a full structural check.</p>
066	<p>NOTICE: msgcnt <i>x</i>: msg 066: V-2-66: DMAPI mount event - buffer</p> <p>■ Description An HSM (Hierarchical Storage Management) agent responded to a DMAPI mount event and returned a message in buffer.</p> <p>■ Action Consult the HSM product documentation for the appropriate response to the message.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
067	<p>WARNING: msgcnt x: mesg 067: V-2-67: mount of <i>device_path</i> requires HSM agent</p> <ul style="list-style-type: none"> ■ Description The file system mount failed because the file system was marked as being under the management of an HSM agent, and no HSM agent was found during the mount. ■ Action Restart the HSM agent and try to mount the file system again.
069	<p>WARNING: msgcnt x: mesg 069: V-2-69: memory usage specified by the vxfs:vxfs_ninode and vxfs:vx_bc_bufhwm parameters exceeds available memory; the system may hang under heavy load</p> <ul style="list-style-type: none"> ■ Description The value of the system tunable parameters—vx_ninode and vx_bc_bufhwm—add up to a value that is more than 66% of the kernel virtual address space or more than 50% of the physical system memory. VxFS inodes require approximately one kilobyte each, so both values can be treated as if they are in units of one kilobyte. ■ Action To avoid a system hang, reduce the value of one or both parameters to less than 50% of physical memory or to 66% of kernel virtual memory. See “Tuning the VxFS file system” on page 44. Note: The tunable parameter vx_ninode is used to set the value of vxfs_ninode.
070	<p>WARNING: msgcnt x: mesg 070: V-2-70: checkpoint <i>checkpoint_name</i> removed from file system <i>mount_point</i></p> <ul style="list-style-type: none"> ■ Description The file system ran out of space while updating a Storage Checkpoint. The Storage Checkpoint was removed to allow the operation to complete. ■ Action Increase the size of the file system. If the file system size cannot be increased, remove files to create sufficient space for new Storage Checkpoints. Monitor capacity of the file system closely to ensure it does not run out of space. See the <code>fsadm_vxfs(1M)</code> manual page.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
071	<p>NOTICE: msgcnt <i>x</i>: msg 071: V-2-71: cleared data I/O error flag in <i>mount_point</i> file system</p> <ul style="list-style-type: none">■ Description The user data I/O error flag was reset when the file system was mounted. This message indicates that a read or write error occurred while the file system was previously mounted. See Message Number 038.■ Action Informational only, no action required.
072	<p>WARNING: msgcnt <i>x</i>: vxfs: msg 072: could not failover for <i>volume_name</i> file system</p> <ul style="list-style-type: none">■ Description This message is specific to the cluster file system. The message indicates a problem in a scenario where a node failure has occurred in the cluster and the newly selected primary node encounters a failure.■ Action Save the system logs and core dump of the node along with the disk image (metasave) and contact your customer support organization. The node can be rebooted to join the cluster.
075	<p>WARNING: msgcnt <i>x</i>: msg 075: V-2-75: replay fsck failed for <i>mount_point</i> file system</p> <ul style="list-style-type: none">■ Description The log replay failed during a failover or while migrating the CFS primary-ship to one of the secondary cluster nodes. The file system was disabled.■ Action Unmount the file system from the cluster. Use <code>fsck</code> to run a full structural check and mount the file system again.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
076	<p>NOTICE: msgcnt <i>x</i>: mesg 076: V-2-76: checkpoint asynchronous operation on <i>mount_point</i> file system still in progress</p> <p>■ Description</p> <p>An EBUSY message was received while trying to unmount a file system. The unmount failure was caused by a pending asynchronous fileset operation, such as a fileset removal or fileset conversion to a nodata Storage Checkpoint.</p> <p>■ Action</p> <p>The operation may take a considerable length of time. You can do a forced unmount, or simply wait for the operation to complete so file system can be unmounted cleanly.</p> <p>See the <code>umount_vxfs(1M)</code> manual page.</p>
077	<p>WARNING: msgcnt <i>x</i>: mesg 077: V-2-77: vx_fshdchange - <i>mount_point</i> file system number fileset, fileset header: checksum failed</p> <p>■ Description</p> <p>Disk corruption was detected while changing fileset headers. This can occur when writing a new inode allocation unit, preventing the allocation of new inodes in the fileset.</p> <p>■ Action</p> <p>Unmount the file system and use <code>fsck</code> to run a full structural check.</p>
078	<p>WARNING: msgcnt <i>x</i>: mesg 078: V-2-78: vx_ilealloc - <i>mount_point</i> file system <i>mount_point</i> fileset (index number) ilist corrupt</p> <p>■ Description</p> <p>The inode list for the fileset was corrupted and the corruption was detected while allocating new inodes. The failed system call returns an ENOSPC error. Any subsequent inode allocations will fail unless a sufficient number of files are removed.</p> <p>■ Action</p> <p>Unmount the file system and use <code>fsck</code> to run a full structural check.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
079	

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
	WARNING: msgcnt x: msg 017: V-2-79: vx_attr_getblk - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_attr_iget - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_attr_indadd - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_attr_indtrunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_attr_iremove - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_bmap - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_bmap_indirect_ext4 - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_delbuf_flush - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_dio_iovec - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_dirbread - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_dircreate - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_dirlook - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_doextop_iau - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_doextop_now - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_do_getpage - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_enter_ext4 - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_exttrunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk
	WARNING: msgcnt x: msg 017: V-2-79: vx_get_alloc - <i>mount_point</i>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
	file system inode <i>inumber</i> marked bad on disk
079 (continued)	<p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_ilsterr - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_indtrunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_iread - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_iremove - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_iremove_attr - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_logwrite_flush - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_oltmount_iget - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_overlay_bmap - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_readnomap - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_reorg_trunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_stablestore - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_tranitimes - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_trunc - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_write_alloc2 - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_write_default - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p> <p>WARNING: msgcnt <i>x</i>: msg 017: V-2-79: vx_zero_alloc - <i>mount_point</i> file system inode <i>inumber</i> marked bad on disk</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
079 (continued)	<ul style="list-style-type: none">■ Description<p>When inode information is no longer dependable, the kernel marks it bad on disk. The most common reason for marking an inode bad is a disk I/O failure. If there is an I/O failure in the inode list, on a directory block, or an indirect address extent, the integrity of the data in the inode, or the data the kernel tried to write to the inode list, is questionable. In these cases, the disk driver prints an error message and one or more inodes are marked bad.</p><p>The kernel also marks an inode bad if it finds a bad extent address, invalid inode fields, or corruption in directory data blocks during a validation check. A validation check failure indicates the file system has been corrupted. This usually occurs because a user or process has written directly to the device or used <i>fsdb</i> to change the file system.</p><p>The <code>VX_FULFSCCK</code> flag is set in the super-block so <i>fsck</i> will do a full structural check the next time it is run.</p>■ Action<p>Check the console log for I/O errors. If the problem is a disk failure, replace the disk. If the problem is not related to an I/O failure, find out how the disk became corrupted. If no user or process is writing to the device, report the problem to your customer support organization. In either case, unmount the file system and use <i>fsck</i> to run a full structural check.</p>
080	<p>WARNING: msgcnt x: vxfs: mesg 080: Disk layout versions older than Version 4 will not be supported in the next release. It is advisable to upgrade to the latest disk layout version now.</p> <p>See the <i>vxupgrade(1M)</i> manual page.</p> <p>See the <i>Veritas Storage Foundation Release Notes</i>.</p> <ul style="list-style-type: none">■ Action<p>Use the <i>vxupgrade</i> command to begin upgrading file systems using older disk layouts to Version 5, then 6, then 7. Consider the following when planning disk layout upgrades:</p>■ Version 2 disk layout file systems have an 8 million inode limit.<p>Images of Version 2 disk layout file systems created by copy utilities, such as <i>dd</i> or <i>volcopy</i>, will become unusable after a disk layout upgrade.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
081	<p>WARNING: msgcnt <i>x</i>: msg 081: V-2-81: possible network partition detected</p> <ul style="list-style-type: none"> ■ Description This message displays when CFS detects a possible network partition and disables the file system locally, that is, on the node where the message appears. ■ Action There are one or more private network links for communication between the nodes in a cluster. At least one link must be active to maintain the integrity of the cluster. If all the links go down, after the last network link is broken, the node can no longer communicate with other nodes in the cluster. Check the network connections. After verifying that the network connections is operating correctly, unmount the disabled file system and mount it again.
082	<p>WARNING: msgcnt <i>x</i>: msg 082: V-2-82: <i>volume_name</i> file system is on shared volume. It may get damaged if cluster is in partitioned state.</p> <ul style="list-style-type: none"> ■ Description If a cluster node is in a partitioned state, and if the file system is on a shared VxVM volume, this volume may become corrupted by accidental access from another node in the cluster. ■ Action These shared disks can also be seen by nodes in a different partition, so they can inadvertently be corrupted. So the second message 082 tells that the device mentioned is on shared volume and damage can happen only if it is a real partition problem. Do not use it on any other node until the file system is unmounted from the mounted nodes.
083	<p>WARNING: msgcnt <i>x</i>: msg 083: V-2-83: <i>mount_point</i> file system log is not compatible with the specified intent log I/O size</p> <ul style="list-style-type: none"> ■ Description Either the specified <code>mount logiosize</code> size is not compatible with the file system layout, or the file system is corrupted. ■ Action Mount the file system again without specifying the <code>logiosize</code> option, or use a <code>logiosize</code> value compatible with the intent log specified when the file system was created. If the error persists, unmount the file system and use <code>fsck</code> to run a full structural check.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
084	<p>WARNING: msgcnt <i>x</i>: msg 084: V-2-84: in <i>volume_name</i> quota on failed during assumption. (stage <i>stage_number</i>)</p> <ul style="list-style-type: none"> ■ Description In a cluster file system, when the primary of the file system fails, a secondary file system is chosen to assume the role of the primary. The assuming node will be able to enforce quotas after becoming the primary. If the new primary is unable to enforce quotas this message will be displayed. ■ Action Issue the quotaon command from any of the nodes that have the file system mounted.
085	<p>WARNING: msgcnt <i>x</i>: msg 085: V-2-85: Checkpoint quota - warning: <i>file_system</i> file system fileset quota hard limit exceeded</p> <ul style="list-style-type: none"> ■ Description The system administrator sets the quotas for Storage Checkpoints in the form of a soft limit and hard limit. This message displays when the hard limit is exceeded. ■ Action Delete Storage Checkpoints or increase the hard limit.
086	<p>WARNING: msgcnt <i>x</i>: msg 086: V-2-86: Checkpoint quota - warning: <i>file_system</i> file system fileset quota soft limit exceeded</p> <ul style="list-style-type: none"> ■ Description The system administrator sets the quotas for Storage Checkpoints in the form of a soft limit and hard limit. This message displays when the soft limit is exceeded. ■ Action Delete Storage Checkpoints or increase the soft limit. This is not a mandatory action, but is recommended.
087	<p>WARNING: msgcnt <i>x</i>: msg 087: V-2-87: vx_dotdot_manipulate: <i>file_system</i> file system <i>inumber</i> inode ddnumber dotdot inode error</p> <ul style="list-style-type: none"> ■ Description When performing an operation that changes an inode entry, if the inode is incorrect, this message will display. ■ Action Run a full file system check using <code>fsck</code> to correct the errors.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
088	<p>WARNING: msgcnt <i>x</i>: msg 088: V-2-88: quota on <i>file_system</i> failed; limits exceed limit</p> <ul style="list-style-type: none"> ■ Description The external quota file, <i>quotas</i>, contains the quota values, which range from 0 up to 2147483647. When quotas are turned on by the <i>quotaon</i> command, this message displays when a user exceeds the quota limit. ■ Action Correct the quota values in the <i>quotas</i> file.
089	<p>WARNING: msgcnt <i>x</i>: msg 089: V-2-89: quota on <i>file_system</i> invalid; disk usage for group/user id <i>uid</i> exceeds sectors sectors</p> <ul style="list-style-type: none"> ■ Description The supported quota limit is up to 2147483647 sectors. When quotas are turned on by the <i>quotaon</i> command, this message displays when a user exceeds the supported quota limit. ■ Action Ask the user to delete files to lower the quota below the limit.
090	<p>WARNING: msgcnt <i>x</i>: msg 090: V-2-90: quota on <i>file_system</i> failed; soft limits greater than hard limits</p> <ul style="list-style-type: none"> ■ Description One or more users or groups has a soft limit set greater than the hard limit, preventing the BSD quota from being turned on. ■ Action Check the soft limit and hard limit for every user and group and confirm that the soft limit is not set greater than the hard limit.
091	<p>WARNING: msgcnt <i>x</i>: msg 091: V-2-91: vx_fcl_truncate - failure to punch hole at offset <i>offset</i> for <i>bytes</i> bytes in File Change Log file; error <i>error_number</i></p> <ul style="list-style-type: none"> ■ Description The vxfs kernel has experienced an error while trying to manage the space consumed by the File Change Log file. Because the space cannot be actively managed at this time, the FCL has been deactivated and has been truncated to 1 file system block, which contains the FCL superblock. ■ Action Re-activate the FCL.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
092	<p>WARNING: msgcnt <i>x</i>: msg 092: V-2-92: vx_mkfcltran - failure to map offset <i>offset</i> in File Change Log file</p> <ul style="list-style-type: none"> ■ Description The vxfs kernel was unable to map actual storage to the next offset in the File Change Log file. This is mostly likely caused by a problem with allocating to the FCL file. Because no new FCL records can be written to the FCL file, the FCL has been deactivated. ■ Action Re-activate the FCL.
095	<p>WARNING: msgcnt <i>x</i>: msg 095: V-2-95: Setting <i>vxfs_ifree_timelag</i> to <i>time</i> since the specified value for <i>vxfs_ifree_timelag</i> is less than the recommended minimum value of <i>time</i>.</p> <ul style="list-style-type: none"> ■ Description The value for <i>vxfs_ifree_timelag</i> specified by the system administrator is less than the recommended minimum value, <i>time</i>, and so the value of <i>vxfs_ifree_timelag</i> has been automatically changed to <i>time</i>. ■ Action No corrective action required on the file system.
096	<p>WARNING: msgcnt <i>x</i>: msg 096: V-2-96: <i>file_system</i> file system fullfsck flag set - <i>function_name</i>.</p> <ul style="list-style-type: none"> ■ Description The next time the file system is mounted, a full <code>fsck</code> must be performed. ■ Action No immediate action required. When the file system is unmounted, run a full file system check using <code>fsck</code> before mounting it again.
097	<p>WARNING: msgcnt <i>x</i>: msg 097: V-2-97: VxFS failed to create new thread (<i>error_number</i>, <i>function_address</i>:<i>argument_address</i>)</p> <ul style="list-style-type: none"> ■ Description VxFS failed to create a kernel thread due to resource constraints, which is often a memory shortage. ■ Action VxFS will retry the thread creation until it succeeds; no immediate action is required. Kernel resources, such as kernel memory, might be overcommitted. If so, reconfigure the system accordingly.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
098	<p>WARNING: msgcnt <i>x</i>: msg 098: V-2-98: VxFS failed to initialize File Change Log for fileset fileset (index number) of <i>mount_point</i> file system</p> <ul style="list-style-type: none">■ Description VxFS mount failed to initialize FCL structures for the current fileset mount. As a result, FCL could not be turned on. The FCL file will have no logging records.■ Action Reactivate the FCL.
099	<p>WARNING: msgcnt <i>x</i>: msg 099: V-2-99: The specified value for <i>vx_ninode</i> is less than the recommended minimum value of <i>min_value</i></p> <ul style="list-style-type: none">■ Description Auto-tuning or the value specified by the system administrator resulted in a value lower than the recommended minimum for the total number of inodes that can be present in the inode cache. VxFS will ignore the newly tuned value and will keep the value specified in the message (VX_MINNINODE).■ Action Informational only; no action required.
100	<p>WARNING: msgcnt <i>x</i>: msg 100: V-2-100: Inode <i>inumber</i> can not be accessed: file size exceeds OS limitations.</p> <ul style="list-style-type: none">■ Description The specified inode's size is larger than the file size limit of the current operating system. The file cannot be opened on the current platform. This can happen when a file is created on one OS and the filesystem is then moved to a machine running an OS with a smaller file size limit.■ Action If the file system is moved to the platform on which the file was created, the file can be accessed from there. It can then be converted to multiple smaller files in a manner appropriate to the application and the file's format, or simply be deleted if it is no longer required.

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
101	<p>WARNING: msgcnt <i>x</i>: msg 101: V-2-101: File Change Log on <i>mount_point</i> for file set <i>index</i> approaching max file size supported. File Change Log will be reactivated when its size hits max file size supported.</p> <p>■ Description</p> <p>The size of the FCL file is approaching the maximum file size supported. This size is platform specific. When the FCL file reaches the maximum file size, the FCL will be deactivated and reactivated. All logging information gathered so far will be lost.</p> <p>■ Action</p> <p>Take any corrective action possible to restrict the loss due to the FCL being deactivated and reactivated.</p>
102	<p>WARNING: msgcnt <i>x</i>: msg 102: V-2-102: File Change Log of <i>mount_point</i> for file set <i>index</i> has been reactivated.</p> <p>■ Description</p> <p>The size of FCL file reached the maximum supported file size and the FCL has been reactivated. All records stored in the FCL file, starting from the current <i>fc_loff</i> up to the maximum file size, have been purged. New records will be recorded in the FCL file starting from offset <i>fs_bsize</i>. The activation time in the FCL is reset to the time of reactivation. The impact is equivalent to File Change Log being deactivated and activated.</p> <p>■ Action</p> <p>Informational only; no action required.</p>
103	<p>WARNING: msgcnt <i>x</i>: msg 103: V-2-103: File Change Log merge on <i>mount_point</i> for file set <i>index</i> failed.</p> <p>■ Description</p> <p>The VxFS kernel has experienced an error while merging internal per-node File Change Log files into the external File Change Log file. Since the File Change Log cannot be maintained correctly without this, the File Change Log has been deactivated.</p> <p>■ Action</p> <p>Re-activate the File Change Log.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
104	<p>WARNING: msgcnt <i>x</i>: msg 104: V-2-104: File System <i>mount_point</i> device <i>volume_name</i> disabled</p> <p>■ Description</p> <p>The volume manager detected that the specified volume has failed, and the volume manager has disabled the volume. No further I/O requests are sent to the disabled volume.</p> <p>■ Action</p> <p>The volume must be repaired.</p>
105	<p>WARNING: msgcnt <i>x</i>: msg 105: V-2-105: File System <i>mount_point</i> device <i>volume_name</i> re-enabled</p> <p>■ Description</p> <p>The volume manager detected that a previously disabled volume is now operational, and the volume manager has re-enabled the volume.</p> <p>■ Action</p> <p>Informational only; no action required.</p>
106	<p>WARNING: msgcnt <i>x</i>: msg 106: V-2-106: File System <i>mount_point</i> device <i>volume_name</i> has BAD label</p> <p>■ Description</p> <p>A file system's label does not match the label that the multi-volume support feature expects the file system to have. The file system's volume is effectively disabled.</p> <p>■ Action</p> <p>If the label is bad because the volume does not match the assigned label, use the <code>vxvset</code> command to fix the label. Otherwise, the label might have been overwritten and the volume's contents may be lost. Call technical support so that the issue can be investigated.</p>
107	<p>WARNING: msgcnt <i>x</i>: msg 107: V-2-107: File System <i>mount_point</i> device <i>volume_name</i> valid label found</p> <p>■ Description</p> <p>The label of a file system that had a bad label was somehow restored. The underlying volume is functional.</p> <p>■ Action</p> <p>Informational only; no action required.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
108	<p>WARNING: msgcnt x: msg 108: V-2-108: vx_dexh_error - <i>error</i>: fileset <i>fileset</i>, directory inode number <i>dir_inumber</i>, bad hash inode <i>hash_inode</i>, seg <i>segment</i> bno <i>block_number</i></p> <p>■ Description</p> <p>The supplemental hash for a directory is corrupt.</p> <p>■ Action</p> <p>If the file system is mounted read/write, the hash for the directory will be automatically removed and recreated. If the removal or recreation fails, subsequent messages indicate the type of problem. If there are no further messages, the removal and recreation of the hash succeeded.</p>
109	<p>WARNING: msgcnt x: msg 109: V-2-109: failed to tune down <i>tunable_name</i> to <i>tunable_value</i> possibly due to <i>tunable_object</i> in use, could free up only up to <i>suggested_tunable_value</i></p> <p>■ Description</p> <p>When the value of a tunable, such as <i>ninode</i> or <i>bufhwm</i>, is modified using the <code>ktune</code> command, sometimes the tunable cannot be tuned down to the specified value because of the current system usage. The minimum value to which the tunable can be tuned is also provided as part of the warning message.</p> <p>■ Action</p> <p>Tune down the tunable to the minimum possible value indicated by the warning message.</p> <p>See “Tuning the VxFS file system” on page 44.</p>
110	<p>WARNING: msgcnt x: msg 110: V-2-110: The specified value for <i>vx_bc_bufhwm</i> is less than the recommended minimum value of <i>recommended_minimum_value</i>.</p> <p>■ Description</p> <p>Setting the <i>vx_bc_bufhwm</i> tunable to restrict the memory used by the VxFS buffer cache to a value that is too low has a degrading effect on the system performance on a wide range of applications. Symantec does not recommend setting <i>vx_bc_bufhwm</i> to a value less than the recommended minimum value, which is provided as part of the warning message.</p> <p>■ Action</p> <p>Tune the <i>vx_bc_bufhwm</i> tunable to a value greater than the recommended minimum indicated by the warning message.</p>

Table B-1 Kernel messages (*continued*)

Message Number	Message and Definition
111	<p>WARNING: msgcnt <i>x</i>: msg 111: V-2-111: You have exceeded the authorized usage (maximum <i>maxfs</i> unique mounted user-data file systems) for this product and are out of compliance with your License Agreement. Please email sales_mail@symantec.com or contact your Symantec sales representative for information on how to obtain additional licenses for this product.</p> <p>■ Description</p> <p>As per your Storage Foundation Basic license agreement, you are allowed to have only a limited number of VxFS file systems, and you have exceeded this number.</p> <p>■ Action</p> <p>Email sales_mail@symantec.com or contact your Symantec sales representative for information on how to obtain additional licenses for this product.</p>
112	<p>WARNING: msgcnt <i>x</i>: msg 112: V-2-112: The new value requires changes to the inode table, which can be made only after a reboot.</p> <p>■ Description</p> <p>When the value of <code>vx_ninode</code> is modified to tune down the number of inodes, the tunable cannot be tuned down in cases that require changes to the VxFS inode tables. The change to the tunable is marked to take effect after the next boot, during which time the change can be committed after editing the inode tables. Tuning down <code>vx_ninode</code> without changes to the VxFS inode table may cause applications to fail with the ENFILE error.</p> <p>■ Action</p> <p>Reboot the system so that the new tunable value can be committed after you make changes to the VxFS inode tables.</p>

About unique message identifiers

VxFS generates diagnostic or error messages for issues not related to the kernel, which are displayed along with a unique message identifier (UMI). Each message has a description and a suggestion on how to handle or correct the underlying problem. The UMI is used to identify the issue should you need to call Technical Support for assistance.

Unique message identifiers

Some commonly encountered UMIs and the associated messages are described on the following table:

Table B-2 Unique message identifiers and messages

Message Number	Message and Definition
20002	<p>UX:vxfs <i>command</i>: ERROR: V-3-20002: <i>message</i></p> <ul style="list-style-type: none">■ Description The command attempted to call <code>stat()</code> on a device path to ensure that the path refers to a character device before opening the device, but the <code>stat()</code> call failed. The error message will include the platform-specific message for the particular error that was encountered, such as "Access denied" or "No such file or directory".■ Action The corrective action depends on the particular error.
20003	<p>UX:vxfs <i>command</i>: ERROR: V-3-20003: <i>message</i></p> <ul style="list-style-type: none">■ Description The command attempted to open a disk device, but the <code>open()</code> call failed. The error message includes the platform-specific message for the particular error that was encountered, such as "Access denied" or "No such file or directory".■ Action The corrective action depends on the particular error.
20005	<p>UX:vxfs <i>command</i>: ERROR: V-3-20005: <i>message</i></p> <ul style="list-style-type: none">■ Description The command attempted to read the superblock from a device, but the <code>read()</code> call failed. The error message will include the platform-specific message for the particular error that was encountered, such as "Access denied" or "No such file or directory".■ Action The corrective action depends on the particular error.
20012	<p>UX:vxfs <i>command</i>: ERROR: V-3-20012: <i>message</i></p> <ul style="list-style-type: none">■ Description The command was invoked on a device that did not contain a valid VxFS file system.■ Action Check that the path specified is what was intended.

Table B-2 Unique message identifiers and messages (continued)

Message Number	Message and Definition
20076	<div>UX:vxfs command: ERROR: V-3-20076: message</div> <div><div>■ Description</div><div>The command called <code>stat()</code> on a file, which is usually a file system mount point, but the call failed.</div><div>■ Action</div><div>Check that the path specified is what was intended and that the user has permission to access that path.</div></div>
21256	<div>UX:vxfs command: ERROR: V-3-21256: message</div> <div><div>■ Description</div><div>The attempt to mount the file system failed because either the request was to mount a particular Storage Checkpoint that does not exist, or the file system is managed by an HSM and the HSM is not running.</div><div>■ Action</div><div>In the first case, use the <code>fscckptadm list</code> command to see which Storage Checkpoints exist and mount the appropriate Storage Checkpoint. In the second case, make sure the HSM is running. If the HSM is not running, start and mount the file system again.</div></div>

Table B-2 Unique message identifiers and messages (*continued*)

Message Number	Message and Definition
21264	<p>UX:vxfs <i>command</i>: ERROR: V-3-21264: <i>message</i></p> <ul style="list-style-type: none">■ Description<p>The attempt to mount a VxFS file system has failed because either the volume being mounted or the directory which is to be the mount point is busy.</p><p>The reason that a VxVM volume could be busy is if the volume is in a shared disk group and the volume is currently being accessed by a VxFS command, such as <code>fsck</code>, on a node in the cluster.</p><p>One reason that the mount point could be busy is if a process has the directory open or has the directory as its current directory.</p><p>Another reason that the mount point could be busy is if the directory is NFS-exported.</p>■ Action<p>For a busy mount point, if a process has the directory open or has the directory as its current directory, use the <code>fuser</code> command to locate the processes and either get them to release their references to the directory or kill the processes. Afterward, attempt to mount the file system again.</p><p>If the directory is NFS-exported, unexport the directory, such as by using the <code>unshare mntpt</code> command on the Solaris operating environment. Afterward, attempt to mount the file system again.</p>
21268	<p>UX:vxfs <i>command</i>: ERROR: V-3-21268: <i>message</i></p> <ul style="list-style-type: none">■ Description<p>This message is printed by two different commands: <code>fsckpt_restore</code> and <code>mount</code>. In both cases, the kernel's attempt to mount the file system failed because of I/O errors or corruption of the VxFS metadata.</p>■ Action<p>Check the console log for I/O errors and fix any problems reported there. Run a full <code>fsck</code>.</p>

Table B-2 Unique message identifiers and messages (*continued*)

Message Number	Message and Definition
21272	<p>UX:vxfs <i>command</i>: ERROR: V-3-21272: <i>message</i></p> <ul style="list-style-type: none"> ■ Description The mount options specified contain mutually-exclusive options, or in the case of a remount, the new mount options differed from the existing mount options in a way that is not allowed to change in a remount. ■ Action Change the requested mount options so that they are all mutually compatible and retry the mount.
23729	<p>UX:vxfs <i>command</i>: ERROR: V-3-23729: <i>message</i></p> <ul style="list-style-type: none"> ■ Description Cluster mounts require the <code>vxfsckd</code> daemon to be running, which is controlled by VCS. ■ Action Check the VCS status to see why this service is not running. After starting the daemon via VCS, try the mount again.
24996	<p>UX:vxfs <i>command</i>: ERROR: V-3-24996: <i>message</i></p> <ul style="list-style-type: none"> ■ Description In some releases of VxFS, before the <code>VxFS mount</code> command attempts to mount a file system, <code>mount</code> tries to read the VxFS superblock to determine the disk layout version of the file system being mounted so that <code>mount</code> can check if that disk layout version is supported by the installed release of VxFS. If the attempt to read the superblock fails for any reason, this message is displayed. This message will usually be preceded by another error message that gives more information as to why the superblock could not be read. ■ Action The corrective action depends on the preceding error, if any.

Disk layout

This appendix includes the following topics:

- [About disk layouts](#)
- [Supported disk layouts and operating systems](#)
- [VxFS Version 4 disk layout](#)
- [VxFS Version 5 disk layout](#)
- [VxFS Version 6 disk layout](#)
- [VxFS Version 7 disk layout](#)

About disk layouts

The disk layout is the way file system information is stored on disk. On VxFS, seven different disk layout versions were created to take advantage of evolving technological developments.

The disk layout versions used on VxFS are:

Version 1	Version 1 disk layout is the original VxFS disk layout provided with pre-2.0 versions of VxFS.	Not Supported
Version 2	Version 2 disk layout supports features such as filesets, dynamic inode allocation, and enhanced security. The Version 2 layout is available with and without quotas support.	Not Supported

Version 3	Version 3 disk layout encompasses all file system structural information in files, rather than at fixed locations on disk, allowing for greater scalability. Version 3 supports files and file systems up to one terabyte in size.	Not Supported
Version 4	Version 4 disk layout encompasses all file system structural information in files, rather than at fixed locations on disk, allowing for greater scalability. Version 4 supports files and file systems up to two terabytes in size.	Supported
Version 5	Version 5 enables the creation of file system sizes up to 32 terabytes. File sizes can be a maximum of 4 billion file system blocks. File systems larger than 2 TB must be created on a Veritas Volume Manager volume. Version 5 also enables setting up to 1024 access control list (ACL) entries.	Supported
Version 6	Version 6 disk layout enables features such as multi-volume support, cross-platform data sharing, named data streams, and File Change Log.	Supported
Version 7	Version 7 disk layout enables support for variable and large size history log records, more than 2048 volumes, large directory hash, and Dynamic Storage Tiering.	Supported

Some of the disk layout versions were not supported on all UNIX operating systems. Version 2 and 3 file systems can still be mounted, but this will be disallowed in future releases. Currently, the Version 4, 5, 6, and 7 disk layouts can be created and mounted. Version 7 is the default disk layout version.

The `vxupgrade` command is provided to upgrade an existing VxFS file system to the Version 4, 5, 6, or 7 disk layout while the file system remains online. You must upgrade in steps from older to newer layouts.

See the `vxupgrade(1M)` manual page.

The `vxfsconvert` command is provided to upgrade Version 2 and 3 disk layouts to the Version 7 disk layout while the file system is not mounted.

See the `vxfsconvert(1M)` manual page.

Supported disk layouts and operating systems

Table C-1 shows which disk layouts supported on the which operating systems.

Table C-1 File system type and operating system versions

Disk Layout		JFS 3.3, HP-UX 11.11	VxFS 3.5, HP-UX 11.11	VxFS 3.5.2, HP-UX 11.23 PI	VxFS 4.1, HP-UX 11.23 PI	VxFS 5.0, HP-UX 11.23 PI	VxFS 5.0, HP-UX 11i v3
Version 1	mkfs	No	No	No	No	No	No
	Local Mount	No	No	No	No	No	No
	Shared Mount	No	No	No	No	No	No
Version 2	mkfs	Yes	No	No	No	No	No
	Local Mount	Yes	Yes	Yes	Yes	No	No
	Shared Mount	No	No	No	No	No	No
Version 3	mkfs	Yes	No	No	No	No	No
	Local Mount	Yes	Yes	Yes	Yes	No	No
	Shared Mount	No	No	No	No	No	No
Version 4	mkfs	Yes	Yes	Yes	Yes	Yes	Yes
	Local Mount	Yes	Yes	Yes	Yes	Yes	Yes
	Shared Mount	No	Yes	No	Yes	No	No
Version 5	mkfs	No	No	Yes	Yes	Yes	Yes
	Local Mount	No	No	Yes	Yes	Yes	Yes
	Shared Mount	No	No	No	Yes	No	No

Table C-1 File system type and operating system versions *(continued)*

Disk Layout		JFS 3.3, HP-UX 11.11	VxFS 3.5, HP-UX 11.11	VxFS 3.5.2, HP-UX 11.23 PI	VxFS 4.1, HP-UX 11.23 PI	VxFS 5.0, HP-UX 11.23 PI	VxFS 5.0, HP-UX 11i v3
Version 6	mkfs	No	No	No	Yes	Yes	Yes
	Local Mount	No	No	No	Yes	Yes	Yes
	Shared Mount	No	No	No	Yes	Yes	Yes
Version 7	mkfs	No	No	No	Yes	Yes	Yes
	Local Mount	No	No	No	Yes	Yes	Yes
	Shared Mount	No	No	No	Yes	Yes	Yes

VxFS Version 4 disk layout

The Version 4 disk layout allows the file system to scale easily to accommodate large files and large file systems.

The original disk layouts divided up the file system space into allocation units. The first AU started part way into the file system which caused potential alignment problems depending on where the first AU started. Each allocation unit also had its own summary, bitmaps, and data blocks. Because this AU structural information was stored at the start of each AU, this also limited the maximum size of an extent that could be allocated. By replacing the allocation unit model of previous versions, the need for alignment of allocation units and the restriction on extent sizes was removed.

The VxFS Version 4 disk layout divides the entire file system space into fixed size allocation units. The first allocation unit starts at block zero and all allocation units are a fixed length of 32K blocks. An exception may be the last AU, which occupies whatever space remains at the end of the file system. Because the first AU starts at block zero instead of part way through the file system as in previous versions, there is no longer a need for explicit AU alignment or padding to be added when creating a file system.

The Version 4 file system also moves away from the model of storing AU structural data at the start of an AU and puts all structural information in files. So expanding

the file system structures simply requires extending the appropriate structural files. This removes the extent size restriction imposed by the previous layouts.

All Version 4 structural files reside in the structural fileset.

The structural files in the Version 4 disk layout are:

File	Description
object location table file	Contains the object location table (OLT). The OLT, which is referenced from the super-block, is used to locate the other structural files.
label file	Encapsulates the super-block and super-block replicas. Although the location of the primary super-block is known, the label file can be used to locate super-block copies if there is structural damage to the file system.
device file	Records device information such as volume length and volume label, and contains pointers to other structural files.
fileset header file	Holds information on a per-fileset basis. This may include the inode of the fileset's inode list file, the maximum number of inodes allowed, an indication of whether the file system supports large files, and the inode number of the quotas file if the fileset supports quotas. When a file system is created, there are two filesets—the structural fileset defines the file system structure, the primary fileset contains user data.
inode list file	Both the primary fileset and the structural fileset have their own set of inodes stored in an inode list file. Only the inodes in the primary fileset are visible to users. When the number of inodes is increased, the kernel increases the size of the inode list file.
inode allocation unit file	Holds the free inode map, extended operations map, and a summary of inode resources.
log file	Maps the block used by the file system intent log.
extent allocation unit state file	Indicates the allocation state of each AU by defining whether each AU is free, allocated as a whole (no bitmaps allocated), or expanded, in which case the bitmaps associated with each AU determine which extents are allocated.
extent allocation unit summary file	Contains the AU summary for each allocation unit, which contains the number of free extents of each size. The summary for an extent is created only when an allocation unit is expanded for use.
free extent map file	Contains the free extent maps for each of the allocation units.

File	Description
quotas files	Contains quota information in records. Each record contains resources allocated either per user or per group.

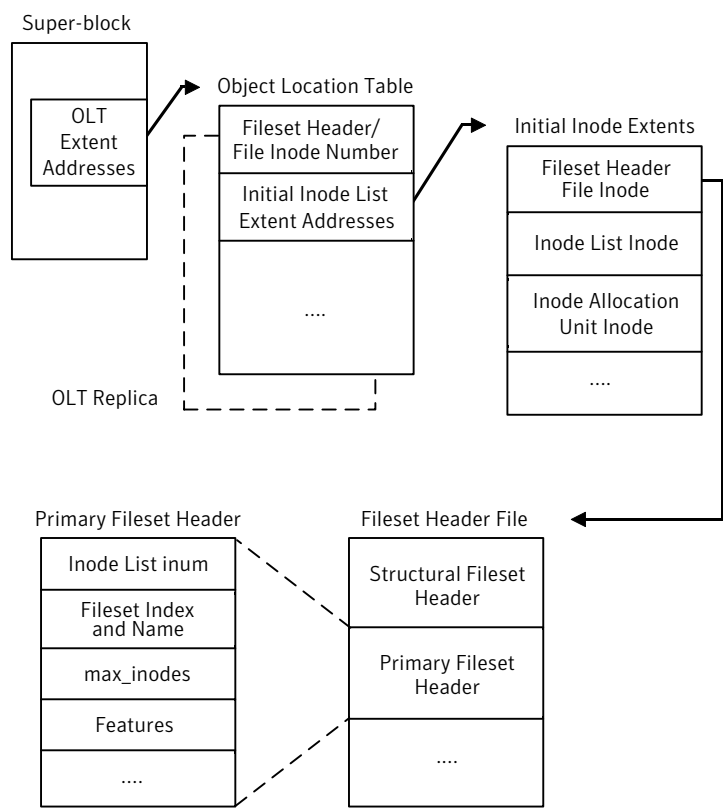
The Version 4 disk layout supports Access Control Lists and Block-Level Incremental (BLI) Backup. BLI Backup is a backup method that stores and retrieves only the data blocks changed since the previous backup, not entire files. This saves times, storage space, and computing resources required to backup large databases.

[Figure C-1](#) shows how the kernel and utilities build information about the structure of the file system.

The super-block location is in a known location from which the OLT can be located. From the OLT, the initial extents of the structural inode list can be located along with the inode number of the fileset header file. The initial inode list extents contain the inode for the fileset header file from which the extents associated with the fileset header file are obtained.

As an example, when mounting the file system, the kernel needs to access the primary fileset in order to access its inode list, inode allocation unit, quotas file and so on. The required information is obtained by accessing the fileset header file from which the kernel can locate the appropriate entry in the file and access the required information.

Figure C-1 VxFS Version 4 disk layout



VxFS Version 5 disk layout

VxFS disk layout Version 5 is similar to Version 4. Structural files in Version 5 are the same in Version 4. However, the Version 5 disk layout supports file systems up to 32 terabytes. For a file system to take advantage of VxFS 32-terabyte support, it must be created on a Veritas Volume Manager volume, and only on a 64-bit kernel operating system. File sizes can be a maximum of 4 billion file system blocks. For 64-bit kernels, the maximum size of the file system you can create depends on the block size:

Block Size	Maximum File System Size
1024 bytes	4,294,967,039 sectors (≈4 TB)

Block Size	Maximum File System Size
2048 bytes	8,589,934,078 sectors (≈8 TB)
4096 bytes	17,179,868,156 sectors (≈16 TB)
8192 bytes	34,359,736,312 sectors (≈32 TB)

If you specify the file system size when creating a file system, the block size defaults to the appropriate value as shown above.

See the `mkfs(1M)` manual page.

See [“About quota files on Veritas File System”](#) on page 110.

The Version 5 disk layout also supports up to 1024 access control list entries.

VxFS Version 6 disk layout

VxFS disk layout Version 6 is similar to Version 5, except that Version 6 enables features such as multi-volume support, cross-platform data sharing, named data streams, and File Change Log. Structural files in Version 6 are the same in Version 5. The Version 6 disk layout can theoretically support files and file systems up to 8 exabytes (2⁶³). For a file system to take advantage of VxFS 8-exabyte support, it must be created on a Veritas Volume Manager volume. For 64-bit kernels, the maximum size of the file system you can create depends on the block size:

Block Size	Currently-Supported Theoretical Maximum File System Size
1024 bytes	68,719,472,624 sectors (≈32 TB)
2048 bytes	137,438,945,248 sectors (≈64 TB)
4096 bytes	274,877,890,496 sectors (≈128 TB)
8192 bytes	549,755,780,992 sectors (≈256 TB)

Note: Sector size in bytes specified by the `DEV_BSIZE` system parameter.

If you specify the file system size when creating a file system, the block size defaults to the appropriate value as shown above.

See the `mkfs(1M)` manual page.

See [“About quota files on Veritas File System”](#) on page 110.

VxFS Version 7 disk layout

VxFS disk layout Version 7 is similar to Version 6, except that Version 7 enables support for variable and large size history log records, more than 2048 volumes, large directory hash, and Dynamic Storage Tiering. The Version 7 disk layout can theoretically support files and file systems up to 8 exabytes (2⁶³). For a file system to take advantage of VxFS 8-exabyte support, it must be created on a Veritas Volume Manager volume. For 64-bit kernels, the maximum size of the file system you can create depends on the block size:

Block Size	Currently-Supported Theoretical Maximum File System Size
1024 bytes	68,719,472,624 sectors (≈32 TB)
2048 bytes	137,438,945,248 sectors (≈64 TB)
4096 bytes	274,877,890,496 sectors (≈128 TB)
8192 bytes	549,755,780,992 sectors (≈256 TB)

Note: Sector size in bytes specified by the `DEV_BSIZE` system parameter.

If you specify the file system size when creating a file system, the block size defaults to the appropriate value as shown above.

See the `mkfs(1M)` manual page.

See [“About quota files on Veritas File System”](#) on page 110.

Glossary

access control list (ACL)	The information that identifies specific users or groups and their access privileges for a particular file or directory.
agent	A process that manages predefined Veritas Cluster Server (VCS) resource types. Agents bring resources online, take resources offline, and monitor resources to report any state changes to VCS. When an agent is started, it obtains configuration information from VCS and periodically monitors the resources and updates VCS with the resource status.
allocation unit	A group of consecutive blocks on a file system that contain resource summaries, free resource maps, and data blocks. Allocation units also contain copies of the super-block.
API	Application Programming Interface.
asynchronous writes	A delayed write in which the data is written to a page in the system's page cache, but is not written to disk before the write returns to the caller. This improves performance, but carries the risk of data loss if the system crashes before the data is flushed to disk.
atomic operation	An operation that either succeeds completely or fails and leaves everything as it was before the operation was started. If the operation succeeds, all aspects of the operation take effect at once and the intermediate states of change are invisible. If any aspect of the operation fails, then the operation aborts without leaving partial changes.
Block-Level Incremental Backup (BLI Backup)	A Symantec backup capability that does not store and retrieve entire files. Instead, only the data blocks that have changed since the previous backup are backed up.
buffered I/O	During a read or write operation, data usually goes through an intermediate kernel buffer before being copied between the user buffer and disk. If the same data is repeatedly read or written, this kernel buffer acts as a cache, which can improve performance. See unbuffered I/O and direct I/O.
contiguous file	A file in which data blocks are physically adjacent on the underlying media.
data block	A block that contains the actual data belonging to files and directories.
data synchronous writes	A form of synchronous I/O that writes the file data to disk before the write returns, but only marks the inode for later update. If the file size changes, the inode will be written before the write returns. In this mode, the file data is guaranteed to be

on the disk before the write returns, but the inode modification times may be lost if the system crashes.

defragmentation	The process of reorganizing data on disk by making file data blocks physically adjacent to reduce access times.
direct extent	An extent that is referenced directly by an inode.
direct I/O	An unbuffered form of I/O that bypasses the kernel's buffering of data. With direct I/O, the file system transfers data directly between the disk and the user-supplied buffer. See buffered I/O and unbuffered I/O.
discovered direct I/O	Discovered Direct I/O behavior is similar to direct I/O and has the same alignment constraints, except writes that allocate storage or extend the file size do not require writing the inode changes before returning to the application.
encapsulation	A process that converts existing partitions on a specified disk to volumes. If any partitions contain file systems, <code>/etc/filesystems</code> entries are modified so that the file systems are mounted on volumes instead. Encapsulation is not applicable on some systems.
extent	A group of contiguous file system data blocks treated as a single unit. An extent is defined by the address of the starting block and a length.
extent attribute	A policy that determines how a file allocates extents.
external quotas file	A <code>quotas</code> file (named <code>quotas</code>) must exist in the root directory of a file system for quota-related commands to work. See <code>quotas</code> file and internal quotas file.
file system block	The fundamental minimum size of allocation in a file system. This is equivalent to the fragment size on some UNIX file systems.
fileset	A collection of files within a file system.
fixed extent size	An extent attribute used to override the default allocation policy of the file system and set all allocations for a file to a specific fixed size.
fragmentation	The on-going process on an active file system in which the file system is spread further and further along the disk, leaving unused gaps or fragments between areas that are in use. This leads to degraded performance because the file system has fewer options when assigning a file to an extent.
GB	Gigabyte (230 bytes or 1024 megabytes).
hard limit	The hard limit is an absolute limit on system resources for individual users for file and data block usage on a file system. See <code>quota</code> .
indirect address extent	An extent that contains references to other extents, as opposed to file data itself. A single indirect address extent references indirect data extents. A double indirect address extent references single indirect address extents.
indirect data extent	An extent that contains file data and is referenced via an indirect address extent.

inode	A unique identifier for each file within a file system that contains the data and metadata associated with that file.
inode allocation unit	A group of consecutive blocks containing inode allocation information for a given fileset. This information is in the form of a resource summary and a free inode map.
intent logging	A method of recording pending changes to the file system structure. These changes are recorded in a circular intent log file.
internal quotas file	VxFS maintains an internal quotas file for its internal usage. The internal quotas file maintains counts of blocks and indices used by each user. See quotas and external quotas file.
K	Kilobyte (210 bytes or 1024 bytes).
large file	A file larger than two one terabyte. VxFS supports files up to 8 exabytes in size.
large file system	A file system larger than one terabytes. VxFS supports file systems up to 8 exabytes in size.
latency	For file systems, this typically refers to the amount of time it takes a given file system operation to return to the user.
metadata	Structural data describing the attributes of files on a disk.
MB	Megabyte (220 bytes or 1024 kilobytes).
mirror	A duplicate copy of a volume and the data therein (in the form of an ordered collection of subdisks). Each mirror is one copy of the volume with which the mirror is associated.
multi-volume file system	A single file system that has been created over multiple volumes, with each volume having its own properties.
MVS	Multi-volume support.
object location table (OLT)	The information needed to locate important file system structural elements. The OLT is written to a fixed location on the underlying media (or disk).
object location table replica	A copy of the OLT in case of data corruption. The OLT replica is written to a fixed location on the underlying media (or disk).
page file	A fixed-size block of virtual address space that can be mapped onto any of the physical addresses available on a system.
preallocation	A method of allowing an application to guarantee that a specified amount of space is available for a file, even if the file system is otherwise out of space.
primary fileset	The files that are visible and accessible to the user.
quotas	Quota limits on system resources for individual users for file and data block usage on a file system. See hard limit and soft limit.

quotas file	The quotas commands read and write the external quotas file to get or change usage limits. When quotas are turned on, the quota limits are copied from the external quotas file to the internal quotas file. See quotas, internal quotas file, and external quotas file.
reservation	An extent attribute used to preallocate space for a file.
root disk group	A special private disk group that always exists on the system. The root disk group is named rootdg.
shared disk group	A disk group in which the disks are shared by multiple hosts (also referred to as a cluster-shareable disk group).
shared volume	A volume that belongs to a shared disk group and is open on more than one node at the same time.
snapshot file system	An exact copy of a mounted file system at a specific point in time. Used to do online backups.
snapped file system	A file system whose exact image has been used to create a snapshot file system.
soft limit	The soft limit is lower than a hard limit. The soft limit can be exceeded for a limited time. There are separate time limits for files and blocks. See hard limit and quotas.
Storage Checkpoint	A facility that provides a consistent and stable view of a file system or database image and keeps track of modified data blocks since the last Storage Checkpoint.
structural fileset	The files that define the structure of the file system. These files are not visible or accessible to the user.
super-block	A block containing critical information about the file system such as the file system type, layout, and size. The VxFS super-block is always located 8192 bytes from the beginning of the file system and is 8192 bytes long.
synchronous writes	A form of synchronous I/O that writes the file data to disk, updates the inode times, and writes the updated inode to disk. When the write returns to the caller, both the data and the inode have been written to disk.
TB	Terabyte (240 bytes or 1024 gigabytes).
transaction	Updates to the file system structure that are grouped together to ensure they are all completed.
throughput	For file systems, this typically refers to the number of I/O operations in a given unit of time.
unbuffered I/O	I/O that bypasses the kernel cache to increase I/O performance. This is similar to direct I/O, except when a file is extended; for direct I/O, the inode is written to disk synchronously, for unbuffered I/O, the inode update is delayed. See buffered I/O and direct I/O.

volume	A virtual disk which represents an addressable range of disk blocks used by applications such as file systems or databases.
volume set	A container for multiple different volumes. Each volume can have its own geometry.
vxfs	The Veritas File System type. Used as a parameter in some commands.
VxFS	Veritas File System.
VxVM	Veritas Volume Manager.

Index

A

- access control lists 24
- allocation policies 64
 - default 64
 - extent 18
 - extent based 18
 - multi-volume support 129

B

- bad block revectoring 37
- blkclear 22
- blkclear mount option 38
- block based architecture 27
- block size 18
- blockmap for a snapshot file system 80
- buffer cache high water mark 46
- buffered file systems 22
- buffered I/O 71

C

- cache advisories 73
- cio
 - Concurrent I/O 43
- closesync 22
- cluster mount 26
- commands
 - cron 30
 - fsadm 29
 - getext 66
 - setext 66
- contiguous reservation 65
- converting a data Storage Checkpoint to a nodata Storage Checkpoint 94
- convosync mount option 35, 39
- copy-on-write technique 83, 87
- creating a multi-volume support file system 126
- creating file systems with large files 42
- creating files with mkfs 190, 192
- cron 30, 49
- cron sample script 50

D

- data copy 70
- data integrity 22
- data Storage Checkpoints definition 88
- data synchronous I/O 38, 71
- data transfer 70
- default
 - allocation policy 64
 - block sizes 18
- default_indir_size tunable parameter 58
- defragmentation 29
 - extent 49
 - scheduling with cron 49
- delaylog mount option 35–36
- device file 259
- direct data transfer 70
- direct I/O 70
- directory reorganization 50
- disabled file system
 - snapshot 81
 - transactions 209
- discovered direct I/O 71
- discovered_direct_iosize tunable parameter 53
- disk layout
 - Version 1 255
 - Version 2 255
 - Version 3 256
 - Version 4 256, 258
 - Version 5 256, 261
 - Version 6 256
 - Version 7 256
- disk space allocation 18
- displaying mounted file systems 196
- Dynamic Storage Tiering
 - multi-volume support 122

E

- edquota
 - how to set up user quotas 206
- encapsulating volumes 123
- enhanced data integrity modes 22

- ENOENT 213
- ENOSPC
 - 101
- ENOTDIR 213
- expansion 30
- extent 18, 63
 - attributes 63
 - indirect 19
 - reorganization 50
- extent allocation 18
 - aligned 64
 - control 63
 - fixed size 63
 - unit state file 259
 - unit summary file 259
- extent size
 - indirect 19
- external quotas file 110

F

- fc_foff 116
- fcl_inode_aging_count tunable parameter 56
- fcl_inode_aging_size tunable parameter 57
- fcl_keeptime tunable parameter 54
- fcl_maxalloc tunable parameter 54
- fcl_winterval tunable parameter 55
- file
 - device 259
 - extent allocation unit state 259
 - extent allocation unit summary 259
 - fileset header 259
 - free extent map 259
 - inode allocation unit 259
 - inode list 259
 - intent log 259
 - label 259
 - object location table 259
 - quotas 260
 - sparse 65
- file change log 54
- file system
 - block size 67
 - buffering 22
 - displaying mounted 196
 - increasing size 198
- fileset
 - header file 259
 - primary 85

- filesystems file 195
- fixed extent size 63
- fixed write size 65
- fragmentation
 - monitoring 49–50
 - reorganization facilities 49
 - reporting 49
- fragmented file system characteristics 49
- free extent map file 259
- free space monitoring 48
- freeze 73
- freezing and thawing, relation to Storage Checkpoints 85
- fsadm 29
 - how to reorganize a file system 200
 - how to resize a file system 198
 - reporting extent fragmentation 49
 - scheduling defragmentation using cron 49
- fsadm_vxfs 43
- fscat 76
- fsck 94
- fsckptadm
 - Storage Checkpoint administration 90
- fstab file
 - editing 195
- fstyp
 - how to determine the file system type 197
- fsvoladm 126

G

- get I/O parameter ioctl 74
- getacl 24
- gettext 66
- global message IDs 210

H

- high water mark 46
- how to access a Storage Checkpoint 92
- how to create a backup file system 202
- how to create a Storage Checkpoint 91
- how to determine the file system type 197
- how to display mounted file systems 195
- how to edit the fstab file 195
- how to edit the vfstab file 195
- how to mount a Storage Checkpoint 92
- how to remove a Storage Checkpoint 92
- how to reorganize a file system 200
- how to resize a file system 198

- how to restore a file system 204
- how to set up user quotas 206
- how to turn off quotas 206
- how to turn on quotas 205
- how to unmount a Storage Checkpoint 94
- how to view quotas 206
- HSM agent error message 234–235
- hsm_write_prealloc 55

I

- I/O
 - direct 70
 - sequential 71
 - synchronous 71
- I/O requests
 - asynchronous 38
 - synchronous 37
- increasing file system size 198
- indirect extent
 - address size 19
 - double 19
 - single 19
- initial_extent_size tunable parameter 56
- inode allocation unit file 259
- inode list error 210
- inode list file 259
- inode table 45
 - internal 45
 - sizes 45
- inodes, block based 18
- intent log 20
 - file 259
 - multi-volume support 123
- Intent Log Resizing 21
- internal inode table 45
- internal quotas file 110
- ioctl interface 63

K

- kctune 45, 47–48
- kernel tunable parameters 44

L

- label file 259
- large files 23, 42
 - creating file systems with 42
 - mounting file systems with 42
- largefiles mount option 42

- local mount 26
- log failure 210
- log mount option 34
- logiosize mount option 37

M

- max_buf_data_size tunable parameter 56
- max_direct_iosize tunable parameter 57
- max_diskq tunable parameter 57
- max_seqio_extent_size tunable parameter 57
- metadata
 - multi-volume support 123
- mincache mount option 35, 38
- mkfs
 - creating files with 190, 192
 - creating large files 43
- modes
 - enhanced data integrity 22
- monitoring fragmentation 49
- mount 21, 43
 - how to display mounted file systems 195
 - how to mount a file system 192
 - mounting a Storage Checkpoint 92
 - pseudo device 93
- mount options 34
 - blkclear 38
 - choosing 34
 - combining 44
 - convosync 35, 39
 - delaylog 23, 35–36
 - extended 21
 - largefiles 42
 - log 22, 34
 - logiosize 37
 - mincache 35, 38
 - nodatainlog 35, 37
 - tmplog 36
- mounted file system
 - displaying 196
- mounting a file system 192
 - option combinations 44
 - with large files 42
- mounting a Storage Checkpoint 94
- mounting a Storage Checkpoint of a cluster file system 94
- msgcnt field 211
- multi-volume support 122
 - creating a MVS file system 126
- multiple block operations 18

N

- name space
 - preserved by Storage Checkpoints 84
- ncheck 120
- nodata Storage Checkpoints 94
- nodata Storage Checkpoints definition 88
- nodatainlog mount option 35, 37

O

- O_SYNC 35
- object location table file 259

P

- parameters
 - default 52
 - tunable 52
 - tuning 51
- performance
 - overall 34
 - snapshot file systems 78
- primary fileset relation to Storage Checkpoints 85
- pseudo device 93

Q

- qio_cache_enable tunable parameter 58
- quota commands 110
- quotaoff
 - how to turn off quotas 206
- quotaon
 - how to turn on quotas 205
- quotas 109
 - exceeding the soft limit 110
 - hard limit 109
 - 107
 - how to view quotas 206
 - soft limit 109
- quotas file 110, 260
- quotas.grp file 110

R

- read-only Storage Checkpoints 92
- read_ahead 59
- read_nstream tunable parameter 53
- read_pref_io tunable parameter 52
- removable Storage Checkpoints definition 89

- reorganization

- directory 50
 - extent 50
- report extent fragmentation 49
- reservation space 63
- Reverse Path Name Lookup 120

S

- sam 45, 47–48
- sequential I/O 71
- setacl 24
- setext 66
- snapof 77
- snapped file systems 24, 75
 - performance 78
 - unmounting 76
- snaread 76
- snapshot 202
 - how to create a backup file system 202
- snapshot file system
 - on CFS 76
- snapshot file systems 24, 75
 - blockmap 80
 - creating 77
 - data block area 80
 - disabled 81
 - errors 224
 - fscat 76
 - fuser 76
 - mounting 77
 - multiple 76
 - performance 78
 - read 76
 - super-block 80
- snapsize 77
- sparse file 65
- storage
 - clearing 38
 - uninitialized 38
- Storage Checkpoints
 - accessing 92
 - administration of 90
 - converting a data Storage Checkpoint to a nodata Storage Checkpoint with multiple Storage Checkpoints 97
 - creating 91
 - data Storage Checkpoints 88
 - definition of 83

Storage Checkpoints *(continued)*

- difference between a data Storage Checkpoint and a nodata Storage Checkpoint 95
 - freezing and thawing a file system 85
 - mounting 92
 - multi-volume support 123
 - nodata Storage Checkpoints 88, 94
 - operation failures
 - 101
 - pseudo device 93
 - read-only Storage Checkpoints 92
 - removable Storage Checkpoints 89
 - removing 92
 - space management
 - 101
 - synchronous vs. asynchronous conversion 94
 - types of 88
 - unmounting 94
 - using the fsck command 94
 - writable Storage Checkpoints 92
- super-block 80
- SVID requirement
- VxFS conformance to 30
- synchronous I/O 71
- system failure recovery 20
- system performance
- overall 34

T

- temporary directories 23
- thaw 74
- Thin Reclamation 27, 50
- tmplog mount option 36
- transaction disabling 209
- tunable I/O parameters 52
 - default_indir_size 58
 - discovered_direct_iosize 53
 - fcl_keeptime 54
 - fcl_maxalloc 54
 - fcl_winterval 55
 - initial_extent_size 56
 - inode_aging_count 56
 - inode_aging_size 57
 - max_buf_data_size 56
 - max_direct_iosize 57
 - max_diskq 57
 - max_seqio_extent_size 57
 - qio_cache_enable 58
 - read_nstream 53

tunable I/O parameters *(continued)*

- read_pref_io 52
 - vx_bc_bufhvm 46
 - vx_maxlink 47
 - vxfs_bc_bufhvm 47
 - write_nstream 53
 - write_pref_io 53
 - write_throttle 60
- tuning I/O parameters 51
- tuning VxFS 44
- typed extents 19

U

- umount command 195
- uninitialized storage, clearing 38
- unmount 94, 210
 - a snapped file system 76

V

- Version 1 disk layout 255
- Version 2 disk layout 255
- Version 3 disk layout 256
- Version 4 disk layout 256, 258
- Version 5 disk layout 256, 261
- Version 6 disk layout 256
- Version 7 disk layout 256
- vfstab file
 - editing 195
- virtual disks 30
- volume sets 124
- VOP_INACTIVE 227
- vx_bc_bufhwm tunable parameter 46
- VX_DSYNC 72
- VX_FREEZE 73
- VX_FULLFCK 210, 212–216, 220–222, 224, 227–228, 231–234, 241
- VX_GETCACHE 73
- vx_maxlink tunable parameter 47
- vx_ninode 45
- VX_SETCACHE 73
- VX_SNAPREAD 76
- VX_THAW 74
- VX_UNBUFFERED 71
- vxdump 66
- VxFS
 - storage allocation 33
- vxfs_bc_bufhwm tunable parameter 47
- vxfs_inotopath 120

- vxfsstat 47
- vxfsu_fcl_sync 55
- vxlsino 120
- vxrestore 66, 204
- vxtunefs
 - changing extent size 19
- vxvset 124

W

- writable Storage Checkpoints 92
- write size 65
- write_nstream tunable parameter 53
- write_pref_io tunable parameter 53
- write_throttle tunable parameter 60