



The XtreemFS Developer Guide

Version 1.0



XtreemFS is developed within the **XtreemOS project**. XtreemOS is a Linux-based Grid operating system that transparently integrates Grid user, VO and resource management traditionally found in Grid Middleware. The XtreemOS project is funded by the European Commission's IST program under contract #FP6-033576.

XtreemFS is available from the **XtreemFS website** (www.XtreemFS.org).

This document is © 2009 by the XtreemOS consortium.

Contents

1	Introduction	1
1.0.1	Document Structure	2
2	XtreemFS Servers	3
2.1	DIR - Directory Service	4
2.2	MRC - Metadata and Replica Catalog	4
	Architecture	4
	Processing Stage	5
	Database Backend	5
	Metadata for Volume Management	6
	Metadata for Files and Directories	7
2.3	OSD - Object Storage Device	15
	Striping	16
	Read-only replication	16
3	Client	17
3.0.1	Architecture	17
	FUSE	18
3.0.2	Implementation	19
	Generated interfaces	20
	Lines of code	20

4	RMS - Replica Management Service	21
4.0.3	Choosing the best replica	21
	Vivaldi Algorithm	22
	Vivaldi in XtremFS	23
	Vivaldi in the OSDs	23
	Vivaldi in clients	24
	Replica Selection with Vivaldi	24
4.0.4	Replica creation	26
	Reactive replica creation with Vivaldi	26
	Proactive replica creation with Oraculo	26
	Integration of Oraculo with OSDs	28
4.0.5	Replica deletion	28
4.0.6	Interaction with the Application Execution Management	28
5	Testing	31
5.0.7	Testing POSIX compliance of XtremFS	31
	How to execute the tests	33
	Results	34
5.0.8	Regression Tests	34
6	Protocol and Interactions	35
6.0.9	Constants	35
6.0.10	Types	36
	Globally Shared Types	36
	struct UserCredentials	37
	struct VivaldiCoordinates	37
	Types Shared between MRC and OSD	37
	struct NewFileSize	37
	struct OSDtoMRCDData	38
	struct OSDWriteResponse	38
	struct StripingPolicy	38

struct Replica	38
struct XLocSet	38
struct XCap	39
struct FileCredentials	39
sequence<FileCredentials> FileCredentialsSet	39
Exceptions	39
exception ProtocolException	40
exception errnoException	40
exception RedirectException	40
exception ConcurrentModificationException	40
exception InvalidArgumentException	40
6.0.11 Directory Service Interface	40
struct AddressMapping	41
sequence<AddressMapping> AddressMappingSet	41
struct Service	41
void xtreemfs_address_mappings_get(string uuid, out AddressMappingSet address_mappings)	41
void xtreemfs_address_mappings_remove(string uuid)	42
uint64_t xtreemfs_address_mappings_set(AddressMappingSet address_mappings)	42
void xtreemfs_checkpoint()	42
uint64_t xtreemfs_global_time_s_get() . .	42
void xtreemfs_service_get_by_type(uint16_t type, out ServiceSet services)	42
void xtreemfs_service_get_by_uuid(string uuid, out ServiceSet services)	43
void xtreemfs_service_get_by_name(string name, out ServiceSet services)	43
uint64_t xtreemfs_service_register(Service service)	43
void xtreemfs_service_deregister(string uuid)	43

void xtreamfs_service_offline(string uuid)	43
void xtreamfs_shutdown()	43
6.0.12 Metadata and Replica Catalog Interface	44
struct Stat	44
struct DirectoryEntry	44
struct StatVFS	44
struct Volume	45
const DEFAULT_ONCRPC_PORT	45
const DEFAULT_ONCRPCS_PORT	45
const DEFAULT_HTTP_PORT	45
exception MRCEException	45
boolean access(string path, uint32_t mode)	45
void chmod(string path, uint32_t mode)	46
void chown(string path, string user_id, string group_id)	46
void create(string path, string user_id, string group_id)	46
void ftruncate(XCap write_xcap, out XCap truncate_xcap)	46
void getattr(string path, out Stat stbuf)	46
void getxattr(string path, string name, out string value)	46
void link(string target_path, string link_path)	47
void listxattr(string path, out StringSet names)	47
mkdir(string path, uint32_t mode)	47
open(string path, uint32_t flags, uint32_t mode, out FileCredentials file_credentials)	47

<code>readdir(string path, out DirectoryEntrySet directory_entries)</code>	47
<code>void removexattr(string path, string name)</code>	48
<code>void rename(string source_path, string target_path, out FileCredentialsSet file_credentials)</code>	48
<code>rmdir(string path)</code>	48
<code>void setattr(string path, Stat stbuf)</code>	48
<code>void setxattr(string path, string name, string value, int flags)</code>	48
<code>void statvfs(string volume_name, out StatVFS stbuf)</code>	48
<code>void symlink(string target_path, string link_path)</code>	49
<code>void unlink(string path, out FileCredentialsSet file_credentials)</code>	49
<code>void utimens(string path, uint64_t atime_ns, uint64_t mtime_ns, uint64_t ctime_ns)</code>	49
<code>void utimens(string path, uint64_t atime_ns, uint64_t mtime_ns, uint64_t ctime_ns)</code>	49
<code>xtreemfs_checkpoint()</code>	50
<code>void xtreemfs_check_file_exists(string volume_id, StringSet file_ids, out string bitmap)</code>	50
<code>void xtreemfs_dump_database(string dump_file)</code>	50
<code>void xtreemfs_get_suitable_osds(string file_id, out StringSet osd_uuids)</code>	50
<code>void xtreemfs_lsvol(out VolumeSet volumes)</code>	50
<code>void xtreemfs_mkvol(Volume volume)</code>	50

void xtreemfs_renew_capability(in XCap old_xcap, out XCap renewed_xcap)	51
void xtreemfs_replica_add(string file_id, Replica new_replica)	51
void xtreemfs_replica_list(string file_id, out ReplicaSet replicas)	51
xtreemfs_replica_remove(string file_id, string osd_uuid, out XCap delete_xcap)	51
xtreemfs_restore_database(string dump_file)	51
xtreemfs_restore_file(string file_path, string file_id, uint64_t file_size, string osd_uuid, int32_t stripe_size)	51
void xtreemfs_rmvol(string volume_name)	52
void xtreemfs_shutdown()	52
void xtreemfs_update_file_size(XCap xcap, OSDWriteResponse osd_write_response)	52
6.0.13 Object Storage Device Interface	52
struct InternalGmax	52
struct ObjectData	52
exception OSDException	53
void read(FileCredentials file_credentials, string file_id, uint64_t object_number, uint64_t object_version, uint32_t offset, uint32_t length, out ObjectData object_data)	53
void truncate(FileCredentials file_credentials, string file_id, uint64_t new_file_size, out OSDWriteResponse osd_write_response)	53
void unlink(FileCredentials file_credentials, string file_id)	53

void write(FileCredentials file_credentials, string file_id, uint64_t object_number, uint64_t object_version, uint32_t offset, uint64_t lease_timeout, ObjectData object_data, out OSDWriteResponse osd_write_response)	54
ObjectData xtreemfs_check_object(FileCredentials file_credentials, string file_id, uint64_t object_number, uint64_t object_version)	54
InternalGmax xtreemfs_internal_get_gmax(FileCredentials file_credentials, string file_id)	54
uint64_t xtreemfs_internal_get_file_size(FileCredentials file_credentials, string file_id)	55
void xtreemfs_internal_truncate(FileCredentials file_credentials, string file_id, uint64_t new_file_size, out OSDWriteResponse osd_write_response)	55
InternalReadLocalResponse xtreemfs_internal_read_local(FileCredentials file_credentials, string file_id, uint64_t object_number, uint64_t object_version, uint64_t offset, uint64_t length)	55
void xtreemfs_cleanup_start(boolean remove_zombies, boolean remove_unavail_volume, boolean lost_and_found)	56
void xtreemfs_cleanup_stop()	56
void xtreemfs_cleanup_status(out string status)	56
void xtreemfs_cleanup_is_running(out boolean is_running)	56
void xtreemfs_cleanup_get_results(out StringSet results)	57
void xtreemfs_cleanup_shutdown()	57
6.0.14 Interactions	57

delete	57
read	57
write	58
fsync	60
mkvol	61
rmvol	61
removeReplica	62

Chapter 1

Introduction

XtreemFS [3] is an object-based [2, 7] file system designed for Grid environments. It is the main distributed file system in the XtreemOS operating system, which relies on XtreemFS for replicated and low-latency file storage between Grid machines.

From a user's perspective, XtreemFS offers a global view on files. Files and directory trees are arranged into volumes. A volume can be mounted at any Grid node where a sufficiently authorized job can access and modify files on the volume. Applications access directories and files on XtreemFS volumes through normal POSIX interfaces (`open`, `read`, etc.) and thus do not require re-compilation in order to work with XtreemFS. This stands in marked contrast with earlier Grid file systems such as GFarm [10], which often forced users to rewrite parts of their applications in order to access files across the Grid via special non-POSIX APIs or to adapt to a non-POSIX file system semantics.

From an administrator's perspective, an XtreemFS installation consists of file system clients running on each user's machine and network-based services for storing and retrieving file metadata and data. The former services are known as Metadata and Replica Services (MRCs), while the latter are called Object Storage Services (OSDs). These services are complemented by the Replica Management Service (RMS), which is responsible for creating replicas on demand in response to changing user access patterns as well as eliminating redundant replicas; and the Object Sharing Service (OSS), which provides transaction-based sharing of volatile memory objects and supports memory-mapped files for XtreemFS.

This deliverable is intended to serve as a developer guide. Its focus is on the current design and implementation of the XtreemFS client and servers, net-

work protocols used between clients and servers, and test suites for XtreamFS.

1.0.1 Document Structure

The report is structured as follows. Sections 2.2, 2, 2.1, 2.3 describe the XtreamFS directory, metadata, and object store services. In section 3 we introduce the new XtreamFS client, which was designed from the ground up to take advantage of the new binary protocol and to remedy numerous performance and scalability problems in the previous revision of the client. Section 4 concerns the XtreamFS Replica Management Service. We conclude with a discussion of recent testing efforts in section 5. Finally, section 6 documents the new binary client-server and server-server protocol, a more efficient and easily-maintained replacement for the text-based protocol of previous releases.

Chapter 2

XtreemFS Servers

All XtreemFS servers (DIR, MRC and OSD) are written in Java and employ an event-based staged design. In addition to the common architecture, they also share the basic libraries like RPC server and client or memory management.

In our design, a stage has one or more threads to do the work. Usually, a stage is used for processing which is blocking (e.g. I/O operations) or consumes larger amounts of CPU time (e.g. checking signatures). Each stage receives requests (events), processes them asynchronously and passes the result to a callback. Operations are the “glue” between the stages. For each client request or internal event, there is an Operation class which implements the logic of the call.

All servers also use a custom method of memory management. To avoid excessive data copying to and from the Java VM, we use direct ByteBuffers which represent raw memory on the heap. These direct ByteBuffers are not managed by the Java garbage collector and excessive allocation and freeing of them causes severe performance problems. To overcome this problem and to reduce overall memory consumption, we use a concurrent BufferPool to allocate ByteBuffers. In addition, we use a wrapper class (called ReusableBuffer) which implements reference counting. It also ensures that ReusableBuffers which have been returned to the pool cannot be used anymore.

The ReusableBuffers must be freed (i.e. returned to the BufferPool) after using them. Failing to do so will cause an error message to be printed on finalization which should help to detect memory leaks. Setting `BufferPool.recordStackTraces` to `true` will add a full stack trace of the allocation to the error message which is useful to locate memory leaks. This option

is only for debugging and should not be used for production due to the performance penalty of recording stack traces on each allocation.

The ONC RPC server and client are used by all three servers as well. Both are implemented using Java's non-blocking network IO NIO and can be used with or without SSL.

The DIR and MRC also use an external key-value store called BabuDB (see <http://babudb.googlecode.com>) to persistently store information. How it is used and how the data is stored in BabuDB is described in the DIR and MRC sections, respectively.

2.1 DIR - Directory Service

The Directory Service (DIR) is the central service registry of XtremFS. All services register and regularly update their registration at the DIR. In addition, it keeps all address mappings which the services need to translate UUIDs to hostname and port. The directory service is also used by the MRCs and OSDs to synchronize their clocks.

Currently, the directory service is a single instance. In the future, this service will be replicated and divided into a hierarchy of DIR services.

Persistent data is stored in BabuDB¹, a non-transactional key-value-store. The service and address mapping records are stored in their XDR representation. This means that the DIR database must be deleted or converted if data structures change.

2.2 MRC - Metadata and Replica Catalog

The Metadata and Replica Catalog (MRC) is responsible for the management of all metadata in an XtremFS installation. Core tasks of the MRC are the management of volumes and directory trees, storage of file and directory metadata and access control enforcement.

Architecture

Aside from the ONC RPC server that listens for incoming client requests, the MRC architecture comprises two core components: the processing stage

¹<http://babudb.googlecode.com>

and the database backend. Each request received by the ONC RPC server is parsed and forwarded to the processing stage, which executes the respective file system logic. Any data that needs to be retrieved or modified during file system logic execution is stored in a database backend.

Processing Stage

The MRC interface consists of multiple so-called *operations*. Each operation relates to an implementation of the logic for the execution of a certain request. There are operations e.g. for opening files, reading directory content, creating volumes, and the like. Operations are named and parametrized similar to their corresponding POSIX calls. To circumvent locking issues in the underlying database, operation execution is serialized for each volume, i.e. no more than one thread may execute operations on a certain volume at the same time.

All operations have a similar composition. First, authorization checks are performed, in order to find out whether the user on behalf of whom the request was sent has sufficient permissions to execute the operation. In case of a positive result, the operation logic is executed. Operation logic execution may involve an arbitrary number of accesses to the underlying database backend. A `readdir` request will e.g. result in a database lookup for the content of a directory, a `setxattr` request will cause an extended attribute of a file to be added in the database.

A detailed description of the interface to the MRC including all operations is given in Sec. 6.0.12.

Database Backend

The database backend is accessed at record level, i.e. at a granularity of single key-value pairs. The creation of a new file could e.g. require several record modifications, since a file metadata object needs to be inserted in the database, a link to the parent directory needs to be established, time stamps of parent directories need to be updated, and so forth. Multiple such records can be combined in an insert group, which causes the insertion of a new set of records to take place in a single step, i.e. atomically.

The database backend implementation is decoupled from the remaining MRC code via an interface, which gives developers the opportunity to implement their own database bindings. The currently used implementation is based on BabuDB. A BabuDB instance may comprise multiple databases, which may

2.2. MRC - METADATA AND REPLICATIONAL XTREEMFS SERVERS

in turn comprise multiple indices. Databases are identified by name strings, whereas indices of a database are serially numbered. Lookups and insertions are directed to single indices of a database; besides normal value lookups for keys, BabuDB supports queries for key prefixes, which provides the basis for an efficient lookup of consecutive key-value pairs.

A range of different indices are used to store XtreamFS metadata. How XtreamFS metadata is mapped to BabuDB indices will be described in the following.

Metadata for Volume Management Volume metadata is stored in a database named V. It is arranged in the following indices:

#	Name	Description
0	Volume ID Index	Maps a volume UUID to a volume metadata entity.
1	Volume Name Index	Maps a volume name to a volume ID.

Volume ID Index		
key		
Element	# Bytes	Description
volumeID	var	the volume ID string
value		
Element	# Bytes	Description
fileAccPolID	2	the file access policy ID for the volume
osdPolID	2	the OSD selection policy ID for the volume
offsVolName	2	the offset position of the 'volName' element, relative to the offset of the buffer's first byte
offsPolArgs	2	the offset position of the 'osdPolArgs' element, relative to the offset of the buffer's first byte
volID	var	the volume's UUID string
volName	var	the volume's name string
osdPolArgs	var	the volume's OSD selection policy argument string

Volume Name Index		
key		
Element	# Bytes	Description
volName	var	the volume name string
value		
Element	# Bytes	Description
volId	var	the volume's UUID string

Metadata for Files and Directories File system metadata of a volume is stored in a BabuDB database with a name equal to the volume's UUID. Various indices are used to manage metadata pertaining to files and directories, which will be described in the following tables. Indices have been designed with the following goals in mind:

- Lookups performed by frequently invoked operations should be as fast as possible, like metadata lookups for a given directory path.
- Database records that are frequently updated should include as little unchanged data as possible.
- Frequently performed database updates should be fast, i.e. involve as little index insertions as possible.
- Indices should contain as little redundancy as possible, in order to minimize database size and memory footprint.

With the aforementioned goals in mind, we decided to have a primary index for the primary metadata of files, which maps a key essentially consisting of a parent directory ID and a file name hash to a value that contains a metadata record. This way, BabuDB prefix lookups for parent directory IDs can be used to efficiently retrieve contents of a directory, while normal lookups can be used to retrieve metadata for a single file. Since POSIX requires support for hard links, i.e. different directory entries pointing to the same metadata, and some operations require a retrieval of file metadata by means of file IDs, we decided to maintain a secondary index that allows a retrieval of metadata by means of a file ID. Other indices are used to store extended attributes and access control lists.

2.2. MRC - METADATA AND REPLICATION METADATA FOR CREEMFS SERVERS

#	Name	Description
0	File Index	Stores primary metadata for a directory entry. Values in the index may be of different kinds: <ul style="list-style-type: none"> • <i>frequently changed metadata</i> - encapsulates all metadata that is frequently modified, such as time stamps or file sizes • <i>rarely changed metadata</i> - encapsulates all metadata that is infrequently changed, such as file names, access modes, or ownership of a file • <i>replica location metadata</i> - encapsulates X-Location lists of files • <i>hard link targets</i> - in case additional hard links exist for one file, the value is a hard link target, i.e. a key in the File ID index. Lookups to file metadata will then be performed in two steps: first, a lookup in the File Index will be performed, in order to retrieve the hard link target; then, metadata will be looked up in the File ID Index.
1	XAttr Index	Contains any extended attributes of files and directories. This includes Softlink targets and default striping policies, since they are mapped to extended attributes.
2	ACL Index	Contains access control list entries of all files.
3	File ID Index	The file ID index is used to retrieve file metadata by means of its ID. If no hard links have been created to a file, the file ID will be mapped to a key in the file index, for which the metadata will have to be retrieved with a second lookup. Such a mapping is necessary for some operations that are based on file IDs instead of path names. If hard links have been created, the file ID will be directly mapped to the three different types of primary file metadata (i.e. rarely and frequently changed metadata, as well as replica locations). In this case, the file's entries in the file index point to the corresponding prefix key in the file ID index.
4	Last ID Index	Contains a single key-value pair that maps a static key to the last file ID that has been assigned to a file. The index ensures that new file IDs are assigned to newly created files or directories.

File Index		
key		
Element	# Bytes	Description
parentID	8	file ID of the parent directory
fileNameHash	4	a hash value of the file name
type	1	type of metadata (0=frequently changed metadata, 1=rarely changed metadata, 2=replica locations, 3=hard link targets)
collCount	2	counter that is incremented with each collision of file name hashes - will be omitted unless multiple file names in the directory have the same hash values
value, type = 0		
Element	# Bytes	Description
fcMetadata	20/12/8	frequently changed metadata associated with the file (see 'fcMetadata' definition), 20 bytes for files & symlinks, 12 for directories, 8 for hard link targets
value, type = 1		
Element	# Bytes	Description
rcMetadata	var	rarely changed metadata associated with the file (see 'rcMetadata' definition)
value, type = 2		
Element	# Bytes	Description
xLocList	var/8	the replica location list associated with the file (see 'xLocList' definition), variable length for files & directories, 8 for hard link targets

2.2. MRC - METADATA AND REPLICATIONALOCREEMFS SERVERS

XAttr Index		
key		
Element	# Bytes	Description
fileID	8	the ID of the file to which the extended attribute has been assigned
ownerHash	4	a hash value of the attribute's owner
attrNameHash	4	a hash value of the attribute name
collCount	2	counter that is incremented with each collision of (ownerHash, attrNameHash) pairs - will be omitted unless different attributes are hashed to the same such pair
value		
Element	# Bytes	Description
offsKey	2	the offset position of the 'attrKey' element, relative to the offset of the buffer's first byte
offsValue	2	the offset position of the 'attrValue' element, relative to the offset of the buffer's first byte
attrOwner	var	the user ID of the attribute's owner
attrKey	var	the attribute key
attrValue	var	the attribute value

ACL Index		
key		
Element	# Bytes	Description
fileID	8	the ID of the file to which the extended attribute has been assigned
entityName	var	the name of the entity associated with the ACL entry
value		
Element	# Bytes	Description
rights	2	the access rights for the entity

File ID Index		
key		
Element	# Bytes	Description
fileID	8	the ID of the file
type	1	type of metadata (0=frequently changed metadata, 1=rarely changed metadata, 2=replica locations, 3=hard link targets)
value, type = 0		
Element	# Bytes	Description
fcMetadata	20/12	frequently changed metadata associated with the file (see 'fcMetadata' definition), 20 bytes for files & symlinks, 12 for directories
value, type = 1		
Element	# Bytes	Description
rcMetadata	var	rarely changed metadata associated with the file (see 'rcMetadata' definition)
value, type = 2		
Element	# Bytes	Description
xLocList	var	the replica location list associated with the file (see 'xLocList' definition)
value, type = 3		
Element	# Bytes	Description
parentID	8	the ID of the parent directory in which the metadata for the file is stored
fileName	var	the file name in the parent directory

2.2. MRC - METADATA AND REPLICATION LOGS

Last ID Index		
key		
Element	# Bytes	Description
'*'	1	the only key in the table
value		
Element	# Bytes	Description
lastFileID	8	the last ID that has been previously assigned to a file or directory

Data types referenced in the index descriptions above are listed in the following:

frequentlyChangedMetadata		
files		
Element	# Bytes	Description
atime	4	file access time stamp in seconds since 1970
ctime	4	file metadata change time stamp in seconds since 1970
mtime	4	file content modification time stamp in seconds since 1970
size	8	file size in bytes
directories		
Element	# Bytes	Description
atime	4	file access time stamp in seconds since 1970
ctime	4	file metadata change time stamp in seconds since 1970
mtime	4	file content modification time stamp in seconds since 1970

rarelyChangedMetadata		
files		
Element	# Bytes	Description
type	1	the type of the entry (0=file, 1=directory)
id	8	file ID
mode	4	POSIX access mode
linkCount	2	number of hard links to the file
w32attrs	8	Win32-specific attributes
epoch	4	current truncate epoch
issEpoch	4	last truncate epoch that has been issued
readOnly	1	a flag indicating whether the file is suitable for read-only replication
offsOwner	2	offset position of the 'owner' element, relative to the offset of the buffer's first byte
offsGroup	2	offset position of the 'group' element, relative to the offset of the buffer's first byte
fileName	var	name of the file
owner	var	user ID of the file's owner
group	var	group ID of the file's owner
directories		
Element	# Bytes	Description
type	1	the type of the entry (0=file, 1=directory)
id	8	file ID
mode	4	POSIX access mode
linkCount	2	number of hard links to the file
w32attrs	8	Win32-specific attributes
offsOwner	2	offset position of the 'owner' element, relative to the offset of the buffer's first byte
offsGroup	2	offset position of the 'group' element, relative to the offset of the buffer's first byte
fileName	var	name of the directory
owner	var	user ID of the directory's owner
group	var	group ID of the directory's owner

2.2. MRC - METADATA AND REPLICATIONAL CREEMFS SERVERS

xLocList		
xLocList		
Element	# Bytes	Description
version	4	version of the X-Locations list
replCount	4	number of replicas in the list
offsUpdPol	4	offset position of the 'updPol' element, relative to the offset of the buffer's first byte
offs1	4	offset positions for all replicas, relative to the offset of the buffer's first byte
...	...	
offsN	4	
xLoc1	var	replicas in the X-Locations list
...	...	
xLocN	var	
updPol	var	update policy string that describes how replica updates will be propagated
xLoc		
Element	# Bytes	Description
offsOsdList	2	offset position of the 'osdList' element, relative to the offset of the buffer's first byte
strPol	var	striping policy associated with the replica
osdList	var	list of all OSDs for the replica
strPol		
Element	# Bytes	Description
stripeSize	4	size of a single stripe (=object) in kB
width	4	number of OSDs for the striping
pattern	var	string containing the striping pattern
osdList		
Element	# Bytes	Description
osdCount	2	number of OSDs in the list
offsOSD1	2	offset positions for all OSD UUIDs, relative to the offset of the buffer's first byte
...	...	
offsOSDn	2	
osdUUID1	var	UUIDs of all OSDs in the list
...	...	
osdUUIDn	var	

2.3 OSD - Object Storage Device

The Object Storage Device (OSD) is responsible for reading and writing objects from/to disk. In addition, it also implements the replication (which is transparent to clients). In this section, we first describe the stages and components of the OSD. We then describe the interaction between OSDs for striped files and for read-only replication.

- **StorageStage and StorageThread**

The StorageStage distributes the request onto a pool of StorageThreads. The allocation of requests is based on the fileID to ensure that all requests for a single file are handled by the same thread. This is necessary to avoid sharing of file metadata across multiple threads.

The StorageThread implements the actual file I/O to access objects on disk. It uses a StorageLayout which is responsible for arranging the objects into files and directories in the underlying file system.

- **PreprocStage**

Analyzes the incoming RPC requests and starts the matching Operation for the requests. It also parses the request arguments (RPC message) based on the Operation. In addition, it parses and validates the signed capability and ensures that the client is authorized to execute the operation. To enhance performance, the PreprocStage keeps a cache of validated Capabilities and XLocation lists. The PreprocStage also keeps a list of open files which is updated whenever a file (i.e. a file's object) is accessed. The list is regularly checked (approx. every mminute) for last access times and files which have timed out will be closed. This close event is sent to the other stages, to allow them to clean their caches. As POSIX requires that a file which is deleted while still opened can still be read or written to, the close event is also used to finally remove data of deleted files.

- **DeletionStage**

This stage is removing the objects on disk for files which have been deleted. This is done directly when the unlink RPC is received or when the file is closed (see PreprocStage).

- **VivaldiStage**

Implements the OSD's Vivaldi component and regularly updates its coordinates. See Sec. 4 for details on Vivaldi.

- **ReplicationStage**
Fetches data from remote OSDs for files which are replicated. For more details about the read-only replication see Sec. 2.3.
- **CleanupThread**
This is not a regular stage, but a background task to scan for orphaned files. If a file is deleted on the MRC but the client fails to delete the file at the OSD, we get so called zombies. To remove them, the OSD has to scan its file system from time to time and check the files at the MRC. How often and when these cleanup operations should be executed depends on the usage pattern of the system (e.g. client's often disconnecting during operations).

Striping

XtreemFS allows files to be striped (distributed) over several OSDs. To ensure correct POSIX semantics in this distributed case, OSDs need to exchange additional information on some write and read operations. We use additional UDP datagrams on write to disseminate file size update hints among OSDs. See [9] for a detailed description of the algorithms used in XtreemFS.

Read-only replication

The read-only replication allows users to replicate their immutable files with very low overhead. Users can set a file to read-only which means that it cannot be modified anymore. This allows users to add replicas on other OSDs which can either be a “full” or a “lazy” replica. For a “full” replica the OSD will automatically fetch all objects of that file. For a “lazy” replica the OSD only fetches the objects when a client tries to read them. Additional prefetching for “lazy” replicas will be added.

Chapter 3

Client

The XtreamFS client connects applications to XtreamFS and acts as a gateway to the XtreamFS directory (DIR), metadata (MRC), and object store (OSD) servers. From a user’s perspective, the client consists of a number of binary programs that reside on the user’s machine. These programs and their functions are summarized in table 3.1.

Command line tool	Function
<code>xtfs_lsvol</code>	list volumes on an XtreamFS MRC server
<code>xtfs_mkvol</code>	create a volume on an XtreamFS MRC server
<code>xtfs_mount</code>	mount an XtreamFS volume
<code>xtfs_rmvol</code>	delete a volume on an XtreamFS MRC server
<code>xtfs_send</code>	send arbitrary RPCs to an XtreamFS server
<code>xtfs_stat</code>	print statistics on an XtreamFS file or directory

Table 3.1: XtreamFS client command line tools

3.0.1 Architecture

The client is structured as a network of message-processing *stages* connected by queues. These stages are similar to those in the XtreamFS servers, and are designed with the same intent: to increase concurrency while avoiding data races (see [11] and previous deliverables for an explanation of stages). Unlike the XtreamFS servers, which is implemented in Java and uses a custom-built set of classes for managing stages, the client is implemented in C++ and relies on a third party platform, Yield ¹ for much of its low-level functionality,

¹<http://yield.googlecode.com/>

including concurrency control in the form of stages as well as platform-specific primitives such as files and sockets.

The stage architecture of the XtreamFS client is depicted in figure 3.1. Note that this particular configuration of stages is specific to `xtfs_mount`, which consists of a set of FUSE entry points and proxy stages for the various XtreamFS servers with which `xtfs_mount` communicates: the directory server (DIRProxy), the metadata server (MRCProxy), and one or more object stores (OSDProxy).

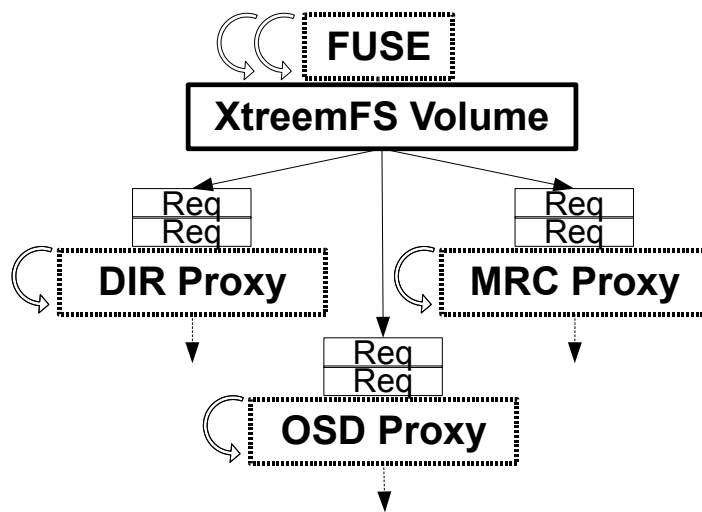


Figure 3.1: Client stages

FUSE

`xtfs_mount` provides a file system interface to applications via FUSE², a library for implementing file systems in userspace. Applications make POSIX system calls such as `open` and `read` into the operating system kernel. A FUSE kernel module translates these calls into messages (i.e. Remote Procedure Calls), which are then passed to a FUSE file system via a pipe. The file system runs as a daemon in an infinite loop, reading and processing messages from the pipe and sending responses back into the kernel. The userspace part of the FUSE library handles most of the nitty gritty details of kernel-userspace communication, so that the file system implementator

²<http://fuse.sourceforge.net/>

can concentrate on implementing file system logic. This is typically done by implementing a set of FUSE callbacks, each of which corresponds to a POSIX system call (and thus a message on the FUSE pipe as well). The FUSE library translates messages to calls on the callbacks supplied by the file system developer and translates return values from the callbacks back into messages for the FUSE kernel module. These FUSE callbacks are the main entry points into `xtfs_mount`.

From a FUSE callback such as `mkdir` the client makes a series of requests (messages) through the various stages shown in figure 3.1. The primary stages in the client are proxies for the servers the `xtfs_mount` instance is connected to: typically a single DIR server and a single MRC server and multiple OSD servers. In the case of `mkdir` the client would send a request to the `MRCProxy` to create the specified directory. Because the FUSE callbacks are synchronous the initial request from a callback to a stage must be synchronous, i.e. the sender must wait for the response before doing any further processing. However, this does not mean that the whole system is synchronous: only the FUSE callback blocks synchronously on a request, allowing the stages in the client to communicate asynchronously (and thus improve performance). Furthermore, since the FUSE callbacks may be multithreaded and reentrant a stage (such as e.g. the `MRCProxy`) can process requests from multiple FUSE callbacks simultaneously. In other words, the concurrency of the client is not limited by the FUSE front end and the number of threads processing FUSE messages.

3.0.2 Implementation

The XtreamFS client is implemented entirely in C++. Aside from the essential components listed above (FUSE callbacks, server proxies) `xtfs_mount` consists of a few support classes such as `Path` (which wraps XtreamFS `volume` file global paths) and XtreamOS integration code such as a pluggable module for retrieving user credentials from an XtreamOS AMS server. Most of this code is shared between the XtreamFS command line tools listed in table 3.1. As mentioned previously, the XtreamFS client also relies heavily on Yield for many low-level classes, such as platform-specific file paths and sockets. The ONC-RPC [8] protocol implementation used to communicate with the XtreamFS servers is also a part of Yield.

Generated interfaces

The synchronous request-response messages exchanges between FUSE callbacks and stages such as the `MRCProxy` are hidden underneath a function call interface. The latter is generated from the same IDL interfaces used by the server. When one of the interface operations is called synchronously on the `MRCProxy` a request is created and filled with the function parameters; the request is sent to the `MRCProxy` stage, where it is processed asynchronously; and the caller blocks waiting for the response, which, when it is received, is unpacked and returned as a normal function return value. The extra level of abstraction allows the FUSE callback interface to be fully agnostic of message sending and receiving, and simply treat the `MRCProxy` as if it were making synchronous remote procedure calls. The FUSE callback for `mkdir` is shown in figure 3.2.

```
bool Volume::mkdir( const YIELD::Path& path, mode_t mode )
{
    mrc_proxy.mkdir( Path( this->name, path ), mode );
    return true;
}
```

Figure 3.2: FUSE callback for `mkdir`

Lines of code

With much of its low-level functionality in `Yield` and other libraries the code base for the `XtreemFS` client is quite minimal, with approximately 2800 lines of hand-written C++ and 4800 lines of C++ automatically generated from the `XtreemFS` IDL interfaces.

Chapter 4

RMS - Replica Management Service

One of the most important mechanisms in XtremFS is the possibility to have several replicas of a file distributed over the Grid. This feature affords data-intensive applications achieving better performance as long as: there is no single access point for the data and mechanisms for parallel access can be exploited. Besides, replication also provides reliability and availability to the filesystem, which is of vital importance for a distributed environment.

However, the usage of Grid resources such as network (where data is transferred across) or storage (where data is stored) are finite, shared, and non-free. Furthermore, the synchronization of the replicas of any given file involves additional overheads, so that mechanisms that keep the tradeoff between the benefits and the extra costs are needed.

For aiming at all of these purposes, we are working on the implementation of the Replica Management Service. This service concerns about: selecting the best replicas for the applications, creating and deleting replicas automatically taking account of how and from where they are accessed and evaluating the maximum number of replicas of any given file.

4.0.3 Choosing the best replica

When a given client (or an OSD) has to access a file, the question is: which replica should it access? It should be able to detect which replica will provide better performance. The idea to solve this problem is to build a virtual 2D space and locate all replicas, OSDs, and clients in it. The distance between two different objects (i.e replica, OSD, or client) is an indicator of the distance

(performance wise) of these two objects. Once a client wants to access a file, it just needs to compute the euclidian distance between itself and all replicas and choose the closer one.

Vivaldi Algorithm

Vivaldi is a light-weight algorithm developed by MIT [1] that allows assigning a position in a coordinate space to every node in a network, so the distance between the coordinates of two nodes predicts the real communication latency between them.

In order to generate a valid coordinate system, it is necessary to determine which space will be used and which formula will be used to calculate the distance between two given points. In our case, it is been proved that implementing a 2-D space, where the Euclidean distance between two coordinates accurately predicts the latency between their corresponding nodes, generates valid results with a really small error probability.

For the algorithm to work correctly, it is also necessary that the nodes of the system keep contacting themselves randomly and indefinitely to re-adjust their position, so any possible change in the network may be reflected. In each re-adjustment, a node contacts a different neighbor, gets its coordinates and modifies its own coordinates, so eventually the Euclidean distance is as similar as possible to the measured round trip time.

On the other hand, once a group of nodes have established a valid coordinate system, it is necessary to use some mechanism that helps to reduce the impact of introducing new nodes, so we avoid them to alter the already stabilized system. That is why Vivaldi keeps in every node, besides the coordinates, a local error that informs about how sure a node is about its position. This way, a node with a steadier position will have a smaller local error and will influence more the rest of nodes when they contact it to readjust their position (figure 4.1).

Once the system is adjusted, any node of the network can determine which nodes are the closest ones with a really simple calculation, in a very short period of time and without generating extra traffic.

Some already developed implementations of Vivaldi can be found in p2psim and in Chord. You might also be interested in Ledlie et al.'s work [6].

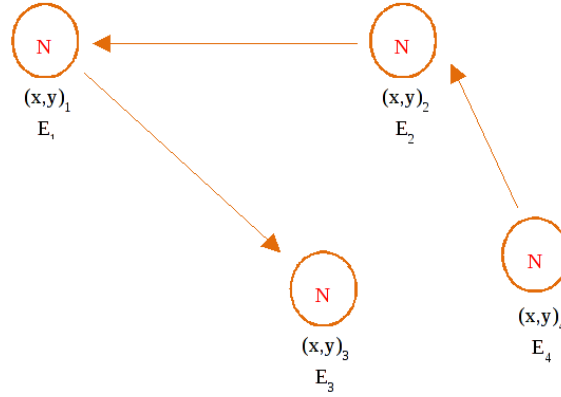


Figure 4.1: Nodes keep recalculating their position

Vivaldi in XtreamFS

As we have different kinds of nodes in our architecture, not all of them work in the same way to integrate Vivaldi. While the clients usually execute during shorter periods of time, the servers are up and running, so the idea is to let the OSDs (at this moment they are the only servers that implement Vivaldi) establish a permanent coordinate system where a client can move through, to find its position.

Vivaldi in the OSDs

An OSD has an independent stage responsible of managing Vivaldi on its side and of providing to the rest of components a couple of valid coordinates that define the position of the node in the current coordinate system.

The stage keeps running indefinitely and periodically contacts a different OSD to ask it for its coordinates and its local error. With that data and the coordinates of the own OSD is possible to compute the Euclidean distance and to compare it with the real RTT measured against the contacted node.

The frequency an OSD readjusts its position is defined by the parameters `MIN_TIMEOUT_RECALCULATE` and `MAX_TIMEOUT_RECALCULATE`. Just after performing a readjustment, the stage typically calculates a random number included in the interval of time defined by those two parameters and sleeps during that number of seconds until the next iteration. This way we try to avoid generating traffic peaks where all the nodes send a request at the same time and to distribute the net use in time.

Larger periods will reduce the overhead in the network but will make the nodes to adjust more slowly to the possible changes in the environment, while smaller ones will require more traffic but will produce a more reactive system.

In each iteration, the introduced stage chooses a node to contact to from a list of available OSDs, which is filled with the information contained in the Directory Service. This list must be updated somehow so the stage can always notice a node going offline.

Vivaldi in clients

In our system, the clients usually execute during a much shorter period of time, so they have to be able to determine their position faster. This can be done because they do not influence the rest of the nodes and they just take some needed info from the already placed OSDs to locate themselves.

In Vivaldi, each node is responsible for its own coordinates and typically has to recalculate them at least a small number of times before they represent the real position in the coordinate system. Even if the set of OSDs is “adjusted”, a client will need to recalculate its position (against one single node each time) several times before having an accurate approximation of its real location. Vivaldi requires that the nodes of the net generate traffic and communicate among themselves.

As in the case of the OSDs, a client also has the parameters `MIN_TIMEOUT_RECALCULATE` and `MAX_TIMEOUT_RECALCULATE` that allow defining the recalculation period. Although the analogue parameters in the OSDs have the same names, they are different parameters and therefore they all must be defined in different files.

Finally, it is important to emphasize that after the first boot of the client, it keeps its coordinates and preserves them among executions, so it remains well located though it mounts and unmounts a lot of different volumes or opens and closes a lot of files. The coordinates are not reinitialized until the client node is rebooted.

Replica Selection with Vivaldi

Until this point we have introduced a mechanism able of establishing a coordinate system where all the nodes of a network have a pair of coordinates that allows them predicting the round trip time to the rest of neighbors.

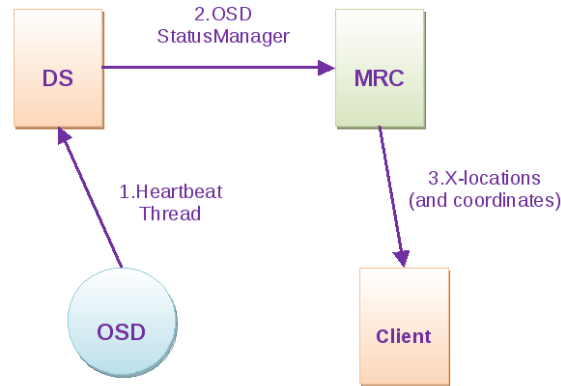


Figure 4.2: Collecting the coordinates

Now it is time to analyze how to take advantage of that information and to describe the current applications of Vivaldi in XtremFS.

Sometimes during the execution of certain operations, the client has to choose which replica access, among several replicas stored in different nodes of the system. The ideal solution proposes to select always the replica that is stored in the closest node, so the accesses can be made within the minimum time. But the problem is that most of the times measuring the RTT against every OSD, for each selection, is not computationally feasible.

Using the coordinates provided by Vivaldi, a client can calculate which replica is the closest one with a practically insignificant delay. At this point the only remaining problem seems to be how to gather all the coordinates so they can be available in the exact moment of the replica selection.

As mentioned earlier, in Vivaldi the coordinates are managed independently by the node they belong to and remain distributed among the different nodes of the network. In order to let the client take advantage of them, it is necessary to collect them in the MRC, so they can be included in every list of x-locations.

In figure 4.2 we show the following process:

1. HeartbeatThread is a component of the OSDs that periodically registers the OSD in the Directory Service. The information that is uploaded each time is determined by the function `getServiceData`, which is defined in the main constructor of the class `OSDRequestDispatcher`.

2. OSDStatusManager is a component of the MRC that regularly queries the Directory Service for available OSDs.
3. During an open operation, the MRC sends to a client a list of x-locations, so it can locate the different replicas associated to a certain file. The x-locations include the corresponding coordinates, so the client can use them to predict the closest replica.

4.0.4 Replica creation

Another important issue regarding replicas is to decide when and where to create a replica. For this functionality we have three different mechanisms. The first one is an explicit request from the user. In this scenario, the RMS will not take any action. The second one is a reactive replica creation. The system will detect that a replica is needed at a given location and will start a replica creation. Finally, in the third case, the system will predict the usage of a file in a location where no replicas are nearby and thus will try to create the replica before it is used. We call to this third mechanism proactive replica creation.

In both cases reactive and proactive, we plan to provide a mechanism able to study the access pattern of files and use it to decide if only a part of the file needs to be replicated (partial replication). This partial replicas will speedup the process of replication because only part of the data will need to be copied to the new location. Nevertheless if we miss-predict the parts of the replica that will be used, we will always be able to populate the missing parts on-demand (done directly by the OSDs).

Reactive replica creation with Vivaldi

In this scenario it is detected when replicas are currently needed in other parts of the Grid. Then, using the distance mechanisms we just described in Section 4.0.3, we will detect if clients request a replica from large distances. So in this case Vivaldi could be used to decide a better location for a replica and create it.

Proactive replica creation with Oraculo

We have implemented a service called Oraculo that carries out data mining on multi-order context models to analyze the file-access patterns and to compute

future accesses. For the implementation of such multi-order context models, Oraculo keeps a trie (or prefix tree) as Kroeger et al. did [4, 5] for centralized environments, where they proved that such structures were effective.

Thus, each file access is modeled as a symbol (i.e. the file path) and it is recorded as a node of the prefix tree with an associated value that represents how many times the chain pattern from root to that node has occurred.

Then, in order to interact with Oraculo, it provides a basic interface for:

1. Adding an event to the trie, given a sequence of the last events seen.
2. Getting a prediction from the trie, given a sequence of the last events seen.

So that when a file access is produced, it can be noticed to Oraculo which computes which parts of the trie must be modified, adding the new event to the corresponding branches or simply increasing the counter of the pattern if it is already present.

Notice that in order to keep the trie scalable Oraculo can prune it in two ways. First, keeping a predefined maximum number of nodes per branch. Thus, whenever a new event goes to a full branch all its nodes divide their corresponding value by two and nodes with a value lower than 1 are deleted from the trie. In the case that no node has been cleaned, the new event is not added. But, obviously the nodes in the branch keep the new value (after the division by two) so in the near future it will be eventually possible to add new events in that branch.

On the other hand, Oraculo also has a maximum of root-branches (also called partitions) to keep the horizontal scalability of the trie. Here we apply a LRU-based algorithm among the partitions taking account of their usage as well.

Finally, Oraculo can predict future file accesses by applying basic data-mining on the trie. It only needs to know some previous accesses to look for patterns based on them. Then, an OSD could eventually use this information to replicate data in advance.

Furthermore, we will propose a decoupled and off-line aggregation of the tries. Once in a while (still to be determined), OSDs could contact other OSDs (only a small subset) to exchange their trie information and build an aggregated one that has the information of all. This mechanism will allow all OSDs to have a more or less global view because what is learned by one OSD will be propagated through several aggregations. We have done some

preliminary tests using this mechanism and seems to work very well with environments of many thousands of nodes.

Regarding the information of the access pattern of files, in most cases the information kept by a single OSD will be enough. Nevertheless, whenever we need the full information of the pattern, we can contact all OSDs that have a replica and aggregate the learned behavior. As we do not expect to keep many replicas of a file, this procedure seems reasonable and scalable.

Integration of Oraculo with OSDs

Unfortunately, the integration of Oraculo with OSDs could not be done yet because we are still evaluating it by using GridSim (a well-known Grid simulator). Once we get significant results with the simulations we will evaluate them and port Oraculo to the OSDs.

4.0.5 Replica deletion

On the one hand, if we want to be able to replicate files whenever needed but still maintain the maximum number for replicas per file, it would be interesting to keep the number of replicas a bit smaller than the maximum. This difference between the maximum and the real number of replicas would allow the system to create replicas whenever needed. On the other hand, if replicas are not used, it would also be nice to have them removed automatically to reduce disk usage in a given node and/or center.

To tackle these two issues we will implement a mechanism that automatically deletes the less important replicas. To know what replicas are less important we will use similar mechanisms than the ones used to create replicas. We will predict future usage using the same kind of tries. In addition we will perform some kind of preventive removal of replicas, which means that whenever a node decides to remove a replica it will inform other OSDs that have it to react accordingly.

4.0.6 Interaction with the Application Execution Management

The last mechanisms that we will implement to manage replicas consists of an interaction with the application execution management system (AEM). This interaction will be done in two steps.

Firstly, AEM analyzes the JSDL of the application and asks to XtreamFS for the locations (coordinates X,Y) of its files' references. Thus AEM computes an optimal coordinate around where the job should be executed. This coordinate is used for AEM to send a request to Resource Selection Service (RSS) for a set of nodes close to it. Of course, RSS also considers other requirements, such as CPU or memory, to decide the resulting set of nodes.

In the second step, the AEM will inform XtreamFS on the final destination of a given job and the files it will use. With this information, the RMS will decide if new replicas need to be created to improve the I/O performance of this job. In addition, and in some cases, it might be that the RMS decides to advance this step from the information obtained in step 1. For instance, this may happen when the list is made of nodes that are close among themselves and one or two replicas could do the job.

Although this mechanism is very good in the sense that no prediction needs to be done, it has a couple of limitations. The first one is that the AEM might not know the files used by a job (it is not a requirement in the job description). The second one is that there might not be enough time from the moment XtreamFS receives the execution location of a job (and the files it uses) and the moment the job starts running. To solve these two cases we have proposed the previous prediction mechanisms (4.0.4).

Chapter 5

Testing

Regular and extensive testing is of vital importance for any file system, in order to improve its reliability, scalability, code quality, stability, performance, etc. Therefore, a new activity inside the WP3.4 has been created, which is focused on all the aspects of testing. The goal of this task is to evaluate functionality and performance of XtreamFS, by exploiting some existing file system test suites.

Nowadays, from a software viewpoint, there are some available tools that are mainly designed for testing local file systems, but there are no ready-made tools available for testing Grid File Systems. Thus, in order to test XtreamFS, either some ad-hoc tools must be developed or some existing distributed testing tools must be extended.

Currently, there are two main guidelines leading the tests performed on XtreamFS: the first one is aimed at evaluating the POSIX compliance of XtreamFS and consequently its functionality; the second one addresses the performance evaluation by stressing XtreamFS with a set of tests and benchmarks (in the following we refer to such kind of tests as regression tests).

In the next subsections we will describe both such activities, the tools used and the results obtained.

5.0.7 Testing POSIX compliance of XtreamFS

One of the goals of XtreamFS is providing a POSIX-compliant filesystem. In order to evaluate the correctness of functionalities, we initially evaluated some available test suites aimed at the POSIX compliance testing. Firstly, we evaluated the "Open Posix Tests Suite" (<https://sourceforge.net/projects/posixtest/>);

it resulted not very suitable for our tests, because it performs only AIO tests but nothing else related to the file system; moreover, we encountered some difficulties during the installation and in particular during the execution of the tests, and for such reasons we discarded it. A second tool that we evaluated was the NTFS-3G suite (<http://www.ntfs-3g.org>), that is a test suite available for the most important operating systems; it includes a POSIX filesystem test environment, the Pawel Jakub Dawidek's POSIX filesystem test suite, that immediately seemed more suitable for our purposes than the first one. For such a reason, we chose to exploit the PJD's POSIX filesystem test suite (PJD-fstest) and to run it over XtremFS. The test suite is available on the Web under a BSD license and we got it from <http://www.ntfs-3g.org/pjd-fstest.html>. PJD-fstest performs almost 3700 regression tests that exhaustively check a wide amount of different scenarios for the following system calls:

- *chmod*: changes the permissions of a files or directories
- *chown*: changes the owner of files or directories
- *link*: creates hard link
- *mkdir*: creates directories
- *mkfifo*: creates fifo named pipes
- *open*: opens a file
- *rename*: renames files or directories
- *rmdir*: removes directories
- *symlink*: creates symbolic links
- *truncate*: truncates files
- *unlink*: removes files

For each system call, the suite contemplates the execution of a set of scripts. Each script performs a set of basic operations, like the creation of a directory, the change of its access rights, the change of its owner, etc., and for each operation it evaluates its return value. If such a value is different than that expected, an error is pointed out. Obviously, the scripts performing the tests for a particular system call are composed of operations targeted for the evaluation of the (hopefully correct) behaviour of that system call. Our

work consists in the automatic execution of the scripts and in the evaluation of the failure events. Then, for each failure, we need to interpret the cause of the problem and reproduce manually the scenario (the sequence of operations) causing it. After this step (that some time hides some difficulties) the problem is pointed out in a bug tracker, in order to be scheduled for a solution.

How to execute the tests

After downloading the tarball of the PJD-fstest, we can simply extract its contents:

```
tar xvzf <package>
```

The most significant contents of the tarball are:

- *fstest.c*: the source code of the main program. It provides an implementation of all the syscalls commented before.
- *Makefile*: the Makefile that we can use to compile the program.
- *tests*: a directory that contains one subdirectory for every syscall. For each of this subdirectories, there is a set of scripts that conform the tests for the corresponding syscall.

Afterwards, we can compile the *fstest.c* file by executing *make*.

To execute the tests, we have implemented a tool that basically automatizes all the process of updating, compiling, and installing XtremFS, running a basic scenario with one Directory Service, one MRC and an OSD, and creating a volume and mounting it on a specific directory.

Once this scenario is up and running, we enter the mount-point where the volume is mounted and execute the tests by using the command *prove* as follows:

```
sudo prove -r <path_to_pjd_suite>/tests
```

The flag *-r* specifies that *prove* should traverse all the directories recursively (so all the tests are executed). Notice that we need root privileges so we also need to edit */etc/fuse.conf* and add the following line:

```
user_allow_other
```

This line allows non-root users to specify the option *allow_other* (or *allow_root*) among the mounting options of fuse.

Results

At the moment, XtreamFS does not support fifos, so mkfifo tests are being ignored.

The tests corresponding to the rest of the syscalls were executed completely and the only errors encountered were due to the lack of implementation of sticky bits on files and directories.

We plan to work on it in next XtreamFS versions, although it is not a very important feature and currently we have to spend our efforts on other more significant issues.

5.0.8 Regression Tests

We use a set of regular file system test tools and custom made tests to automatically check the XtreamFS development version (the svn trunk) every night. This test environment can also be used to manually run these tests.

The main test script is `trunk/bin/xtfs_test` and can be executed to run all tests automatically or to start/stop a test environment.

To start a test environment with all XtreamFS servers and clients, run

```
> trunk/bin/xtfs_test --start
```

This script will put all data and logfiles in the current working directory.

After setting up the test environment, the tests in `trunk/tests/` can be executed individually by calling the test scripts in the mounted XtreamFS volume.

```
> python trunk/tests/01_simple_metadata.sh
```

To shutdown all servers and unmount the clients after testing, execute

```
> trunk/bin/xtfs_test --stop
```

To clean up all data and logfiles use

```
> trunk/bin/xtfs_test --clean
```

If you want to run all tests automatically, run the test script in auto mode

```
> trunk/bin/xtfs_test --autotest
```

Chapter 6

Protocol and Interactions

XtreemFS uses ONC RPC[8] for executing remote operations. Interfaces and records are defined in a subset of CORBA IDL. Yidl¹ is used to generate the code for the interfaces and records in C++ and Java.

To build the Java classes from the interfaces:

1. export PYTHONPATH to point to your yidl source directory:
`export PYTHONPATH=/home/user/yidl/src`
2. execute `bin/generate_xtreemfs_java.py`

6.0.9 Constants

Globally shared constants are defined in `interfaces/constants.idl`.

`ACCESS_CONTROL_POLICY_NULL` don't use any access policy (on the MRC).

This will allow all users to do everything on the volume.

`ACCESS_CONTROL_POLICY_POSIX` use standard POSIX permissions (user, group, others) on the volume.

`ACCESS_CONTROL_POLICY_VOLUME` similar to POSIX permissions but the permission for the root (/) is used for the entire volume.

`ACCESS_CONTROL_POLICY_DEFAULT` the policy to use in e.g. `mkvol` if nothing is specified.

¹<http://code.google.com/p/yidl/>

ONCRPC_SCHEME scheme for URLs.

ONCRPCS_SCHEME scheme for URLs when using SSL.

ONCRPC_AUTH_FLAVOR constant to use for ONC RPC auth_flavor to indicate XtremFS auth. If present, a UserCredentials record is sent in auth_opaque

OSD_SELECTION_POLICY_SIMPLE only OSDs which are alive and which have more than 2GB free space are used.

OSD_SELECTION_POLICY_DEFAULT the policy to use in e.g. mkvol if nothing is specified.

REPL_UPDATE_PC_NONE no replication is used

REPL_UPDATE_PC_RDONLY read-only replication

SERVICE_TYPE_MRC for DIR service registry, service is an MRC

SERVICE_TYPE OSD for DIR service registry, service is an OSD

SERVICE_TYPE_VOLUME for DIR service registry, service is a volume

STRIPING_POLICY_RAID0 RAID0 (striping)

STRIPING_POLICY_DEFAULT the policy to use in e.g. mkvol if nothing is specified.

STRIPING_POLICY_STRIPE_SIZE_DEFAULT default stripe size in KB to use if nothing is specified.

STRIPING_POLICY_WIDTH_DEFAULT default striping width (number of OSDs) to use if nothing is specified.

SYSTEM_V_FCNTL_H_0... POSIX constants

6.0.10 Types

Globally Shared Types

Globally shared data structures are defined in `interfaces/types.idl`.

struct UserCredentials

User information sent in the ONC RPC `opaque_auth` body if XtremFS authentication is used. How the `userID` and `groupIDs` look like depends on the policy used in the client which translates the local `uid/gid`.

<code>user_id</code>	globally unique <code>userID</code>
<code>group_ids</code>	list of globally unique <code>groupIDs</code> (must contain at least one entry)
<code>password</code>	admin password (in cleartext) required for some operations (e.g. <code>mkvol</code>)

struct VivaldiCoordinates

Structure used to exchange Vivaldi coordinates between components, also used in UDP packets for measuring latency between XtremFS clients and OSDs.

<code>x_coordinate</code>	x coordinate
<code>y_coordinate</code>	y coordinate
<code>local_error</code>	confidence in correctness of x/y coordinates

Types Shared between MRC and OSD

Types that are mainly shared between MRC and OSD are defined in `interfaces/mrc_osd_types.idl`.

struct NewFileSize

Sent by the OSD in response to a file modification operation if the file size has changed. A client may cache these updates and send them to the MRC when renewing a capability, on `fsync/flush` and `close`. The client needs only to send the most recent record it received from the OSD for a given file. Most recent means that: `(size_in_bytes' > size_in_bytes AND truncate_epoch' == truncate_epoch)` OR `(truncate_epoch' > truncate_epoch)`

The client should update its local file size cache with the `NewFileSize` records received from the OSD. The client should use the *locally cached* file size on `stat` rather than the result from the MRC to ensure that local processes see their own modifications.

<code>size_in_bytes</code>	the new file size in bytes
<code>truncate_epoch</code>	truncate epoch in which this operation was executed (used by the MRC for ordering updates)

struct OSDtoMRCData

Data sent by the OSD to the client which is expected to pass it on to the MRC. When the data should be passed to the MRC depends on the `cached_policy`. This feature is currently not used.

<code>cached_policy</code>	describes how the client is allowed to cache the data (when to send it to the MRC)
<code>data</code>	opaque data

struct OSDWriteResponse

Record containing file size updates and/or OSDtoMRCData. Returned from all data-modifying operations.

<code>new_file_size</code>	contains no record or at most one record if the file size changed
<code>opaque_data</code>	contains 0 or more records

struct StripingPolicy

Describes how a replica (one copy of the file) is split into objects.

<code>policy</code>	describes the scheme to use for distributing the objects among the OSDs, e.g. RAID0 for simple round robin striping.
<code>stripe_size</code>	the size of the objects in kilobytes, must be ≥ 4
<code>width</code>	the number of OSDs to use for striping, must be ≥ 1 .

struct Replica

Describes a single copy of a file.

<code>striping_policy</code>	the striping policy to use for this replica.
<code>replication_flags</code>	value depends on the replication policy, e.g. to indicate a full or lazy replica.
<code>osd_uuids</code>	ordered (!) list of OSDs holding objects of the file.

struct XLocSet

Describes a complete file together with all replicas (copies) and how they are kept consistent.

<code>replicas</code>	list of the file's replicas (i.e. list of Replica structs)
<code>version</code>	incremented by the MRC on each modification of the list. Used by the OSD to reject clients working with outdated lists.
<code>repUpdatePolicy</code>	the policy used for keeping replicas in sync.
<code>read_only_file_size</code>	the size of the file in bytes, used only for read-only replication.

struct XCap

Security token which is issued by the MRC and authorizes a client to execute operations on a file at the OSDs.

<code>file_id</code>	file for which the capability can be used
<code>access_mode</code>	POSIX access mode for which client is authorized (e.g. read only, delete, write, truncate).
<code>expires_s</code>	absolute timestamp when the capability becomes invalid (seconds since epoch).
<code>client_identity</code>	the client identity set by the MRC, currently the client's IP address.
<code>truncate_epoch</code>	the file's current truncate epoch.
<code>server_signature</code>	the MRC's signature for the capability which is used by the OSD to validate the XCap. Signature is created using shared secret specified in the MRC and OSD configuration.

struct FileCredentials

A record containing the XLocSet and XCap for a file. Required for most OSD operations.

<code>xlocs</code>	the XLocSet
<code>xcap</code>	the capability

sequence<FileCredentials> FileCredentialsSet

Used by the MRC to return no or at most one FileCredentials record.

Exceptions

Exceptions that may be thrown in connection with an RPC are defined in `interfaces/exceptions.idl`.

exception ProtocolException

Thrown on ONC RPC errors (e.g. GARBAGE_ARGS)

accept_stat ONC RPC accept_stat value
error_code POSIX errno, if available
stack_trace optional, for debugging only

exception errnoException

Thrown by the MRC to indicate a POSIX error.

error_code POSIX errno, if available
error_message optional text message
stack_trace optional, for debugging only

exception RedirectException

Thrown by the DIR, MRC and OSD to redirect the client to another service. Use e.g. for master slave replication to direct the client to the current master.

to_uuid service to contact

exception ConcurrentModificationException

Thrown by the DIR if a record was modified by another service on the meantime.

stack_trace optional, for debugging only

exception InvalidArgumentException

Thrown by the DIR if an input value is not acceptable.

error_message error message describing the correct values.

6.0.11 Directory Service Interface

The Directory Service interface is defined in `interfaces/dir_interface.idl`.

`struct AddressMapping`

Maps a service UUID to protocol, hostname/IP and port. A service can have multiple mappings for different networks (e.g. inside a cluster with private IP addresses). At the moment only “*” is supported for `match_network` which indicates a match for all networks.

<code>uuid</code>	the service UUID
<code>version</code>	the record’s version, used by the DIR to detect concurrent modifications
<code>protocol</code>	the protocol used by the service
<code>address</code>	resolvable hostname or IP address in text form
<code>port</code>	port on which the service listens
<code>match_network</code>	for future use, must be *
<code>ttl_s</code>	time to live in seconds, indicates how long this record can be cached before it is re-fetched from the DIR

`sequence<AddressMapping> AddressMappingSet`

Future releases of XtreamFS will support multi-network setups to ease the usage of XtreamFS in shared public/private network environments often found in clusters.

`struct Service`

Information on a service registered at the DIR.

<code>uuid</code>	the service UUID
<code>version</code>	the record’s version, used by the DIR to detect concurrent modifications
<code>type</code>	service type (see 6.0.9)
<code>name</code>	human readable name of the service; for volumes: the unique volume name
<code>last_updated_s</code>	timestamp of the last time (in seconds since epoch) the service updated its entry at the DIR. Used as a coarse-grained heartbeat-signal.
<code>data</code>	a map of additional data which depends on the service (e.g. MRC of a volume or free space of an OSD)

```
void xtreamfs_address_mappings_get( string uuid, out AddressMappingSet  
address_mappings )
```

Get an address mapping for the service specified by `uuid`.

`uuid` the service UUID
`out address_mappings` empty, if no mapping exists, one (or more) records otherwise

`void xtreemfs_address_mappings_remove(string uuid)`

Remove an address mapping from the DIR.

`uuid` the service UUID

`uint64_t xtreemfs_address_mappings_set(AddressMappingSet address_mappings)`

Updates the address mappings for a service.

`address_mappings` the new mappings. The UUID in all records must be the same. The version must be 0 for a new mapping or the version obtained with the last read from the DIR.

returns the new version of the mapping

throws `ConcurrentModificationException` if the record was updated (version incremented by DIR) between reading and updating the record.

`void xtreemfs_checkpoint()`

Forces the DIR to create a BabuDB checkpoint. This operation does not block, the checkpoint is created asynchronously. The admin password must be sent via the XtreamFS authentication.

`uint64_t xtreemfs_global_time_s_get()`

Returns the current system time on the DIR in seconds since epoch. Used to synchronize MRCs and OSDs to the global XtreamFS system time. The DIR system should be synchronized with a precise clock using e.g. `ntp`.

returns system time in seconds since UNIX epoch.

`void xtreemfs_service_get_by_type(uint16_t type, out ServiceSet services)`

Get all services of a specific type registered at the DIR.

`type` the service type to return

`out services` all matching services

`void xtreemfs_service_get_by_uuid(string uuid, out ServiceSet services)`

Get the service information for a service with a specific UUID.

`uuid` the service uuid
`out services` one record, if the service is registered, empty list otherwise

`void xtreemfs_service_get_by_name(string name, out ServiceSet services)`

Get the service information for a service with a specific name.

`name` the service's name
`out services` one record, if the service is registered, empty list otherwise

`uint64_t xtreemfs_service_register(Service service)`

Update a service registration at the DIR. Updates the `last_update_s` field of the service.

`service` the service's data. The UUID must be the service's UUID, the version must be 0 for a new service or the version obtained with the last read from the DIR.
returns the new version of the mapping
throws `ConcurrentModificationException` if the record was updated (version incremented by DIR) between reading and updating the record.

`void xtreemfs_service_deregister(string uuid)`

Removes the service registry entry for the service from the DIR.

`uuid` the service uuid

`void xtreemfs_service_offline(string uuid)`

Sets the `last_update_s` field to 0 which indicates that the service was taken offline.

`uuid` the service uuid

`void xtreemfs_shutdown()`

Shuts down the DIR service, does not force a checkpoint of the database. The admin password must be sent via the XtreamFS authentication.

6.0.12 Metadata and Replica Catalog Interface

The MRC interface is defined in `interfaces/mrc_interface.idl`.

struct Stat

Contains information about a file, directory or symbolic link that is sent to the client in response to a `getattr` request.

<code>mode</code>	the file's current access mode
<code>nlink</code>	the number of hard links to the file
<code>uid</code>	the numeric UID of the file's owner (just for compatibility reasons, will not be filled)
<code>gid</code>	the numeric GID of the file's owner (just for compatibility reasons, will not be filled)
<code>unused_dev</code>	(just for compatibility reasons, will not be filled)
<code>size</code>	the current file size in bytes
<code>atime_ns</code>	the file's atime in nanos
<code>mtime_ns</code>	the file's mtime in nanos
<code>ctime_ns</code>	the file's ctime in nanos
<code>user_id</code>	the XtreamFS user ID string of the file's owner
<code>group_id</code>	the XtreamFS group ID string of the file's owner
<code>file_id</code>	the XtreamFS file ID
<code>link_target</code>	the target path for symbolic links
<code>truncate_epoch</code>	the file's current truncate epoch
<code>attributes</code>	a set of Win32 specific file attributes

struct DirectoryEntry

Contains information about a directory entry that is sent to the client in response to a `readdir` request.

<code>name</code>	the name of the directory entry
<code>stbuf</code>	a buffer of type <code>struct Stat</code> that contains information about the file

struct StatVFS

Contains information about a mounted XtreamFS volume, which is sent to the client in response to a `statvfs` request.

bsize	the file system's block size (1024)
bfree	the number of free blocks
fsid	the file system ID (volume ID)
namelen	maximum file name length (1024)

struct Volume

Contains information about a volume.

name	the volume name
mode	the access mode for the volume's parent directory
osd_selection_policy	the ID of the OSD selection policy for the volume
default_striping_policy	the ID of the default striping policy for the volume
id	the volume UUID
owner_user_id	the XtreamFS user ID of the volume's owner
owner_group_id	the XtreamFS group ID of the volume's owner

const DEFAULT_ONCRPC_PORT

Constant defining the default MRC ONC RPC port.

const DEFAULT_ONCRPCS_PORT

Constant defining the default MRC ONC RPC port for SSL.

const DEFAULT_HTTP_PORT

Constant defining the default MRC HTTP port.

exception MRCException

Thrown by all MRC operations.

boolean access(string path, uint32_t mode)

Checks access to a file or directory. Responds with **true** if access is granted, **false**, otherwise.

path	the path to the file or directory
mode	the access flags to check

```
void chmod( string path, uint32_t mode )
```

Changes the access mode of a file or directory.

`path` the path to the file or directory
`mode` the new access mode

```
void chown( string path, string user_id, string group_id )
```

Changes the owner of a file or directory.

`path` the path to the file or directory
`user_id` the new owner ID
`group_id` the new owning group ID

```
void create( string path, string user_id, string group_id )
```

Creates a new file.

`path` the path to the new file
`mode` the initial access mode for the new file

```
void ftruncate( XCap write_xcap, out XCap truncate_xcap )
```

Issues a new truncate capability for an open file.

`write_xcap` a valid Capability with write permissions to the file
`out truncate_xcap` a new capability with write and truncate permissions, which has to be used for subsequent operations

```
void getattr( string path, out Stat stbuf )
```

Returns information on a file or directory.

`path` the path to the file or directory
`out stbuf` a buffer containing information on the file or directory

```
void getxattr( string path, string name, out string value )
```

Returns the value of an extended attribute of a file or directory.

`path` the path to the file or directory
`name` the name of the attribute
`out value` the attribute value


```
void link( string target_path, string link_path )
```

Creates a new hard link to an existing file.

target_path the path to the existing file
link_path the path defining where the new hard link shall be created

```
void listxattr( string path, out StringSet names )
```

Returns the set of extended attributes assigned to a file or directory.

path the path to the file or directory
names the list of attribute names assigned to the file or directory

```
mkdir( string path, uint32_t mode )
```

Creates a new directory.

path the path to the new directory
mode the initial access mode for the new directory

```
open( string path, uint32_t flags, uint32_t mode, out FileCredentials  
file_credentials )
```

Opens a file by performing an access check and issuing a new Capability for OSD access in case of success.

path the path to the file
flags a set of flags specifying the kind of access that
 is requested
mode initial access mode for a newly created file in
 case **flags** contains **O_CREAT**
out file_credentials a set of file credentials containing the file's X-
 Locations list and the newly issued Capability

```
readdir( string path, out DirectoryEntrySet directory_entries )
```

Lists the content of a directory, including all file metadata.

path the path to the directory
out directory_entries a list containing all nested directory entries for
 the given directory

```
void removexattr( string path, string name )
```

Removes an extended attribute from a file or directory.

`path` the path to the file or directory
`name` the name of the attribute to remove

```
void rename( string source_path, string target_path, out FileCredentialsSet  
file_credentials )
```

Renames a path.

`source_path` the former path to the file or directory
`target_path` the new path to the file or directory
`out file_credentials` contains an X-Locations list and deletion Capability in case the target path was overwritten

```
rmkdir( string path )
```

Removes an empty directory.

`path` the path to the directory

```
void setattr( string path, Stat stbuf )
```

Sets metadata of a file or directory. Currently, this call is only used to set the file's Win32 attributes.

`path` the path to the file or directory
`stbuf` a buffer containing the metadata to set

```
void setxattr( string path, string name, string value, int flags )
```

Sets an extended attribute of a file or directory.

`path` the path to the file or directory
`name` the name of the attribute
`value` the new value for the attribute
`flags` a set of system flags associated with the attribute (currently ignored)

```
void statvfs( string volume_name, out StatVFS stbuf )
```

Returns information about a volume.

`volume_name` the volume name
`stbuf` a buffer containing information about the volume

```
void symlink( string target_path, string link_path )
```

Creates a symbolic link.

`target_path` the target path
`link_path` the path for the new symbolic link

```
void unlink( string path, out FileCredentialsSet file_credentials )
```

Unlinks a file from a directory. If no more links to the file exist, file metadata will be deleted.

`path` the path to the file
`out file_credentials` a set of file credentials containing a deletion Capability and X-Locations list in case file metadata needs to be deleted.

```
void utimens( string path, uint64_t atime_ns, uint64_t mtime_ns,  
uint64_t ctime_ns )
```

Sets the POSIX time stamps of a file or directory.

`path` the path to the file or directory
`atime_ns` the new access time stamp in nanos (will be ignored if set to 0)
`mtime_ns` the new modification time stamp in nanos (will be ignored if set to 0)
`ctime_ns` the new change time stamp in nanos (will be ignored if set to 0)

```
void utimens( string path, uint64_t atime_ns, uint64_t mtime_ns,  
uint64_t ctime_ns )
```

Sets the POSIX time stamps of a file or directory.

`path` the path to the file or directory
`atime_ns` the new access time stamp in nanos (will be ignored if set to 0)
`mtime_ns` the new modification time stamp in nanos (will be ignored if set to 0)
`ctime_ns` the new change time stamp in nanos (will be ignored if set to 0)

`xtreemfs_checkpoint()`

Enforces the creation of a new database checkpoint. The call blocks until checkpoint creation has completed.

`void xtreemfs_check_file_exists(string volume_id, StringSet file_ids, out string bitmap)`

Checks for a set of file IDs whether the given files exist in the given volume. This call is necessary for cleanup purposes.

`volume_id` the volume's UUID
`file_ids` a list of file IDs to check
`out bitmap` a string containing '1's for each file in `file_ids` that exists, and '0's for each file that does not exist, in the same order as the file IDs given in `file_ids`

`void xtreemfs_dump_database(string dump_file)`

Creates an XML dump of the MRC database on the server.

`dump_file` the path at which to store XML dump file on the MRC's local file system

`void xtreemfs_get_suitable_osds(string file_id, out StringSet osd_uuids)`

Returns a list of suitable OSDs for the given file. The call can be used to find OSDs that are suitable for new replicas of a file.

`file_id` the file ID
`out osd_uuids` a list of OSD UUIDs that can be used for the file

`void xtreemfs_lsvol(out VolumeSet volumes)`

Returns a list of all volumes stored on the MRC.

`out volumes` a list containing information on each volume stored on the MRC

`void xtreemfs_mkvol(Volume volume)`

Creates a new volume.

`volume` information about the volume to create

```
void xtreemfs_renew_capability( in XCap old_xcap, out XCap renewed_xcap )
```

Extends the validity of a capability. The capability must be valid in order to be renewed.

old_xcap the capability to be renewed
renewed_xcap a new capability with the same properties except for an
 extended validity period

```
void xtreemfs_replica_add( string file_id, Replica new_replica )
```

Adds a new replica to a file. The file's read-only flag must be set to `true`.

file_id the file ID
new_replica the new replica to be added

```
void xtreemfs_replica_list( string file_id, out ReplicaSet replicas )
```

Returns the list of replicas of a file. Information on each of the replicas in the list includes the striping policy and X-Loc list.

file_id the file ID
out replicas the list of replicas

```
xtreemfs_replica_remove( string file_id, string osd_uuid, out XCap  
delete_xcap )
```

Removes a replica from a file.

file_id the file ID
osd_uuid the UUID of the head OSD
delete_xcap a capability for deleting the data associated with the
 replica on the OSD

```
xtreemfs_restore_database( string dump_file )
```

Restores the MRC database from an XML dump. When the operation is invoked, no volumes may exist in the current database.

dump_file the path to the XML dump file on the MRC's local file system

```
xtreemfs_restore_file( string file_path, string file_id, uint64_t  
file_size, string osd_uuid, int32_t stripe_size )
```

Restores a file from the given metadata.

<code>dump_file_path</code>	the path associated with the restored file
<code>file_id</code>	the ID associated with the restored file
<code>file_size</code>	the size associated with the restored file
<code>osd_uuid</code>	the OSD on which the file content is stored
<code>stripe_size</code>	the stripe size associated with the restored file

`void xtreemfs_rmvol(string volume_name)`

Deletes a volume, including the metadata of all nested files and directories.

`volume_name` the name of the volume to delete

`void xtreemfs_shutdown()`

Gracefully terminates the MRC with all its sub-components.

`void xtreemfs_update_file_size(XCap xcap, OSDWriteResponse osd_write_response)`

Updates the size of a file in response to an OSD write operation.

`osd_write_response` the response from the OSD, which may contain a new file size

6.0.13 Object Storage Device Interface

The OSD interface is defined in `interfaces/osd_interface.idl`.

`struct InternalGmax`

Sent by OSDs to determine the actual file size of a striped file.

<code>epoch</code>	the file's latest truncate epoch the OSD knows
<code>last_object_id</code>	last object number known by the OSD
<code>file_size</code>	locally known file size in bytes

`struct ObjectData`

Sent by OSDs to determine the actual file size of a striped file.

<code>data</code>	object data (file content)
<code>checksum</code>	checksum for data
<code>zero_padding</code>	zeros to append to data (padding objects, POSIX sparse file semantics)
<code>invalid_checksum_on_osd</code>	true if the OSD detected corrupted on-disk data

exception `OSDException`

Thrown by all OSD operations.

<code>error_code</code>	see class <code>org.xtreemfs.osd.ErrorCodes</code> for a list of error codes
<code>error_message</code>	optional, human readable error message
<code>stack_trace</code>	optional, for debugging only

```
void read( FileCredentials file_credentials, string file_id, uint64_t
object_number, uint64_t object_version, uint32_t offset, uint32_t
length, out ObjectData object_data )
```

Reads on object.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id
<code>object_number</code>	the object requested (first object is 0)
<code>object_version</code>	for future use
<code>offset</code>	offset within object
<code>length</code>	number of bytes to read (offset+length must be < object size)
<code>out object_data</code>	the object data read from disk. If less data (data + zero padding) is returned, this indicates an EOF.

```
void truncate( FileCredentials file_credentials, string file_id,
uint64_t new_file_size, out OSDWriteResponse osd_write_response )
```

Truncates a file to the specified length. The client must have a capability valid for truncating. For files which are striped over more than one OSD, this operation must be executed at the *head OSD* which is the first OSD in the replica's OSD list.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id
<code>new_file_size</code>	the new size of the file in bytes to which the file is truncated
<code>out osd_write_response</code>	information which should be passed to the MRC.

```
void unlink( FileCredentials file_credentials, string file_id )
```

Deletes all objects of the file. If the file is currently open (in use) the objects will be deleted on close. This operation returns immediately, the objects

are deleted by the OSD asynchronously. The client must have a capability valid for deleting. For files which are striped over more than one OSD, this operation must be executed at the *head OSD* which is the first OSD in the replica's OSD list.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id

```
void write( FileCredentials file_credentials, string file_id, uint64_t
object_number, uint64_t object_version, uint32_t offset, uint64_t
lease_timeout, ObjectData object_data, out OSDWriteResponse osd_write_respons
```

Writes an object.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id
<code>object_number</code>	the object to write (first object is 0)
<code>object_version</code>	for future use
<code>offset</code>	offset within object
<code>lease_timeout</code>	for future use (timestamp of client lease timeout in seconds since epoch)
<code>object_data</code>	the object data to write into the object; zero_padding is ignored.

```
ObjectData xtreamfs_check_object( FileCredentials file_credentials,
string file_id, uint64_t object_number, uint64_t object_version )
```

Similar to read. The OSD reads the object and validates the checksum but doesn't send the actual data. Used by the file system scrubber to check data integrity.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id
<code>object_number</code>	the object requested (first object is 0)
<code>object_version</code>	for future use
<code>returns</code>	the size of the object in <code>zero_padding</code> and whether the data is corrupted

```
InternalGmax xtreamfs_internal_get_gmax( FileCredentials file_credentials,
string file_id )
```

Returns the locally known truncate epoch, number of objects and file size for a file. Used by an OSD to determine the file size for a file which is striped over more than one OSD.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id
returns	OSD-local file size information

```
uint64_t xtreemfs_internal_get_file_size( FileCredentials file_credentials,
string file_id )
```

Returns the actual file size on the OSD(s). For files which are striped over more than one OSD, this operation must be executed at the *head OSD* which is the first OSD in the replica's OSD list. Used by file system scrubber tools. Should be used to update file size before marking a file read-only.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id
returns	the actual file size of the file

```
void xtreemfs_internal_truncate( FileCredentials file_credentials,
string file_id, uint64_t new_file_size,
out OSDWriteResponse osd_write_response )
```

Used by the head OSD to truncate the file on all OSDs for a file which is striped across more than one OSD. May only be used by OSDs.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id
<code>new_file_size</code>	the new size of the file in bytes to which the file is truncated
<code>out osd_write_response</code>	information which should be passed to the MRC.

```
InternalReadLocalResponse xtreemfs_internal_read_local( FileCredentials
file_credentials, string file_id, uint64_t object_number, uint64_t
object_version, uint64_t offset, uint64_t length )
```

Used by OSDs to fetch objects from other OSDs for replicated files. May only be used by OSDs. This method does not follow POSIX semantics by add padding data but sends the raw object data from disk.

<code>file_credentials</code>	XLocSet and Capability for the file
<code>file_id</code>	the file's file Id
<code>object_number</code>	the object requested (first object is 0)
<code>object_version</code>	for future use
<code>offset</code>	offset within object
<code>length</code>	number of bytes to read
<code>returns</code>	raw object data from disk

```
void xtreamfs_cleanup_start(boolean remove_zombies,  
boolean remove_unavail_volume, boolean lost_and_found )
```

Starts the cleanup process on an OSD. The cleanup process will check for all files on the OSD's disk if they still exist in the MRC. Requires an admin password in the XtreamFS authentication data.

<code>remove_zombies</code>	delete files which have been deleted on the MRC
<code>remove_unavail_volume</code>	delete files if the MRC holding the volume is not available (DANGEROUS!)
<code>lost_and_found</code>	do not delete files but re-create them in a lost+found directory

```
void xtreamfs_cleanup_stop()
```

Aborts the OSD cleanup process. Requires an admin password in the XtreamFS authentication data.

```
void xtreamfs_cleanup_status( out string status )
```

Returns a human readable status string from the cleanup process. Requires an admin password in the XtreamFS authentication data.

<code>out status</code>	human readable status text (in English)
-------------------------	---

```
void xtreamfs_cleanup_is_running( out boolean is_running )
```

Check if the cleanup process is running. Requires an admin password in the XtreamFS authentication data.

<code>out is_running</code>	true, if the process is running
-----------------------------	---------------------------------

```
void xtreamfs_cleanup_get_results( out StringSet results )
```

Returns a list of messages from the cleanup process. Requires an admin password in the XtreamFS authentication data.

```
    out results    list of messages
```

```
void xtreamfs_cleanup_shutdown()
```

Shuts down the OSD. Requires an admin password in the XtreamFS authentication data.

6.0.14 Interactions

This section illustrates the interactions between XtreamFS clients and servers.

delete

Files are deleted as described in the following (see Fig. 6.1):

1. The client receives a delete request from the VFS. It removes the file on the MRC via the `unlink` operation and receives the file credentials, which contain the globally unique XtreamFS fileID, a deletion capability and the replica locations list.
2. The client initiates the deletion of file content by invoking the `unlink` operation on the head OSD (i.e. the first OSD of a stripe).
3. The head OSD delays the deletion until all clients have closed the file, i.e. all capabilities known to the head OSD have timed out. In turn, the head OSD initiates the deletion of file objects on the remaining OSDs via `unlink`.

read

Files are read as described in the following (see Fig. 6.2):

1. The client receives an `open` request from the VFS. It opens the file on the MRC for reading, and receives the file credentials, which contain the globally unique XtreamFS fileID, a read capability and the replica locations list.

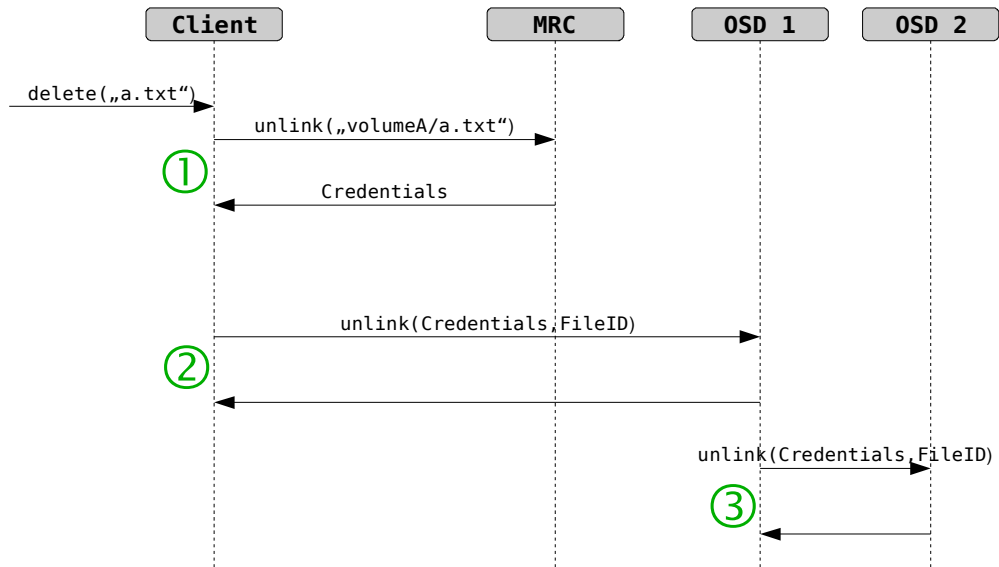


Figure 6.1: Deleting a file

2. The client sends a request for reading Y bytes of data from the offset Z of the file identified by the fileID. The OSD returns a buffer containing object data, as well as additional information like checksum failure notifications or padding flags.
3. If multiple read requests are send, the client has to ensure that the capability is renewed before it times out, in order to keep the file open.
4. The client reads more data from the file, e.g. another object.
5. The client receives a `close` call. There is no need to explicitly close the file on the servers; this is implicitly done when the capabilities in the OSD cache time out.

write

Files are written as described in the following (see Fig. 6.3):

1. The client receives an `open` request from the VFS. It opens the file on the MRC for writing, and receives the file credentials, which contain the globally unique XtreamFS fileID, a write capability and the replica locations list.

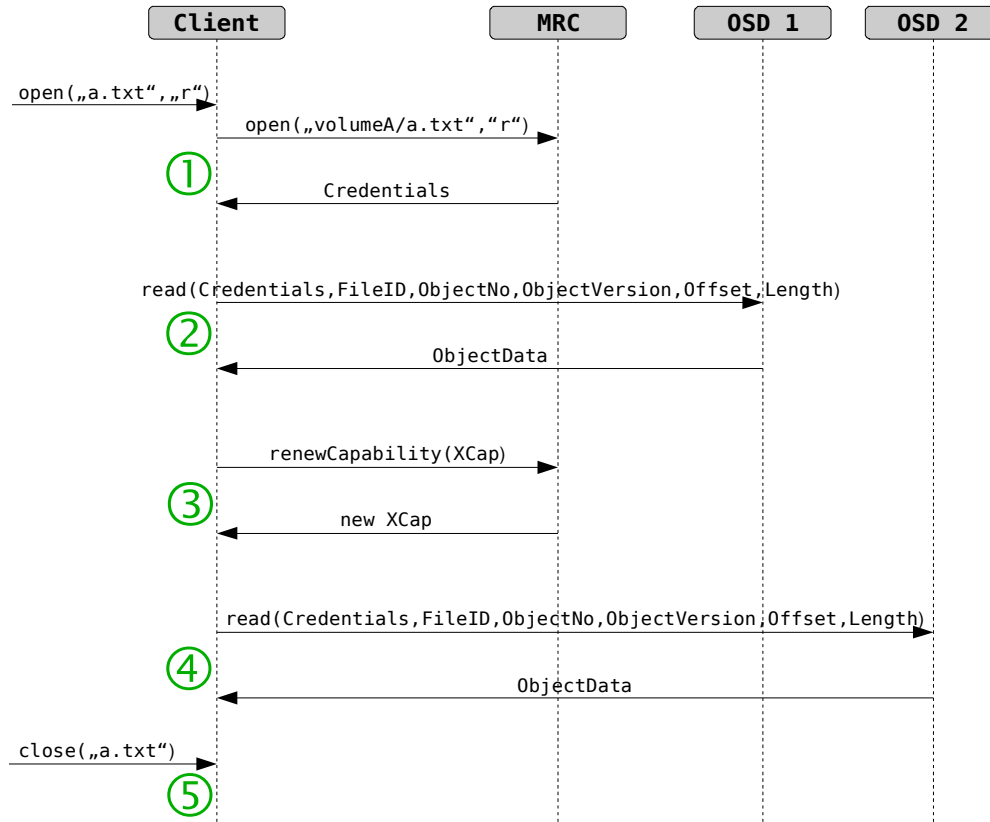


Figure 6.2: Reading a file

2. The client sends a buffer containing the data to the OSD on which the object is stored. In return, the OSD sends an `OSDWriteResponse`, which must be cached and sent to the MRC when the file is closed or `fsync` is called. The client can also decide to send pending filesize updates from time to time between capability renewals, as the lifetime of a capability can be in the range of tens of minutes.
3. If multiple write requests are send, the client has to ensure that the capability is renewed before it times out, in order to keep the file open.
4. The client appends data to the file by sending more `write` requests, each being answered with an `OSDWriteResponse`.
5. The client receives a `close` call. It sends any pending `OSDWriteResponses` to the MRC, in order to update the file size. There is no need to explicitly close the file on the servers; this is implicitly done when the capabilities in the OSD cache time out.



Figure 6.3: Writing a file

fsync

Files are fsync'ed as described in the following (see Fig. 6.4):

1. The file was opened and modified.
2. The client receives a **fsync** request from the VFS. If the file is opened all pending data will be written to the OSD.
3. The client sends any pending **OSDWriteResponses** to the MRC, in order to update the file size.
4. The file is further modified or closed.

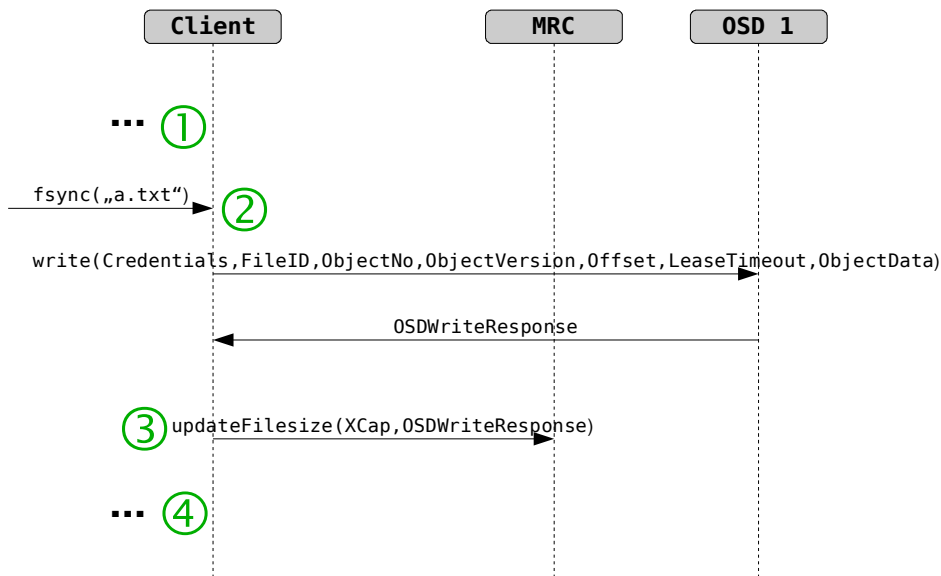


Figure 6.4: Synchronizing data with the underlying device

mkvol

New volumes are created as described in the following (see Fig. 6.5):

1. The client receives an `mkvol` request. In response, it creates the volume on the MRC.
2. The MRC registers the volume at the DIR.
3. If no problems occur, the MRC responds to the client with an acknowledgment.

rmvol

Volumes are deleted as described in the following (see Fig. 6.6):

1. The client receives a `rmvol` request from the console. In response, it removes the volume on the MRC.
2. The MRC deregisters the volume at the DIR.
3. If no problems occur, the MRC responds to the client with an acknowledgment.

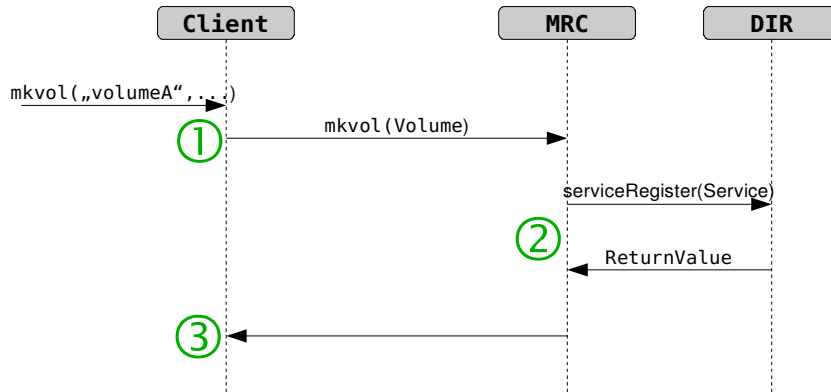


Figure 6.5: Creating a new volume

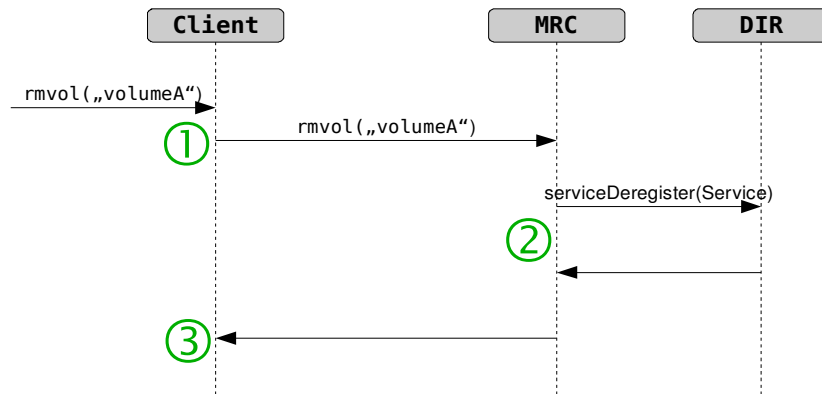


Figure 6.6: Deleting an existing volume

removeReplica

Single replicas of a file are deleted as described in the following (see Fig. 6.7):

1. The client receives a **removeReplica** request from the console. It sends a request to the MRC to remove the replica with a matching head OSD. In response, it receives the file credentials, which contain the globally unique XtreamFS fileID, a deletion capability and the list of replica locations.
2. The client initiates the deletion of the replica's file content by invoking the **unlink** operation on the head OSD (i.e. the first OSD of a stripe).
3. The head OSD delays the deletion until all clients have closed the file, i.e. all capabilities known to the head OSD have timed out. In turn,

the head OSD initiates the deletion of file objects on the remaining OSDs via `unlink`.

4. If a client finds out that its replica locations list for the file is outdated, it has to retrieve the new replica locations list from the MRC.

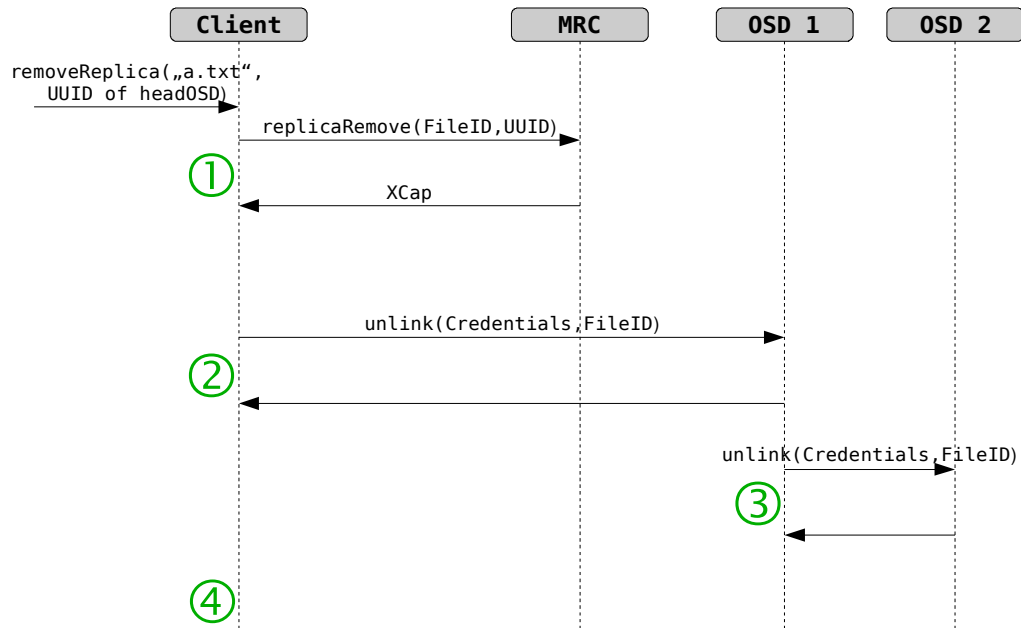


Figure 6.7: Deleting a replica of a file

Index

DIR, [3](#), [4](#), [36](#), [40–43](#), [61](#)

FUSE, [18–20](#)

IDL, [20](#), [35](#)

MRC, [1](#), [3–5](#), [16](#), [17](#), [19](#), [35–42](#), [44](#),
[45](#), [50–53](#), [55–63](#)

ONC RPC, [4](#), [5](#), [35–37](#), [40](#), [45](#)

OSD, [1](#), [3](#), [14–17](#), [19](#), [36–42](#), [45](#), [47](#),
[50–60](#), [62](#), [63](#)

OSS, [1](#)

POSIX, [1](#), [5](#), [7](#), [13](#), [15](#), [16](#), [18](#), [19](#), [35](#),
[36](#), [39](#), [40](#), [49](#), [52](#), [55](#)

UDP, [16](#), [37](#)

Bibliography

- [1] F. Dabek, R. Cox, F. Kaashoek, and R. Morris. Vivaldi: A decentralized network coordinate system. *ACM SIGCOMM Computer Communication Review*, 34(4):15–26, 2004.
- [2] Michael Factor, Kalman Meth, Dalit Naor, Ohad Rodeh, and Julian Satran. Object storage: The future building block for storage systems. In *2nd International IEEE Symposium on Mass Storage Systems and Technologies*, 2005.
- [3] Felix Hupfeld, Toni Cortes, Bjoern Kolbeck, Jan Stender, Erich Focht, Matthias Hess, Jesus Malo, Jonathan Marti, and Eugenio Cesario. XtremFS: a case for object-based storage in Grid data management. In *3rd VLDB Workshop on Data Management in Grids, co-located with VLDB 2007*, 2007.
- [4] T.M. Kroeger. *Predicting File System Actions From Reference Patterns*. PhD thesis, University of California, 1996.
- [5] T.M. Kroeger and D.D.E. Long. Design and Implementation of a Predictive File Prefetching Algorithm. *Proceedings of the General Track: 2002 USENIX Annual Technical Conference table of contents*, pages 105–118, 2002.
- [6] J. Ledlie, P. Gardner, and M. Seltzer. Network coordinates in the wild. In *Proc. of NSDI*, 2007.
- [7] M. Mesnier, G. Ganger, and E. Riedel. Object-based storage. *IEEE Communications Magazine*, 8:84–90, 2003.
- [8] R. Srinivasan. Rpc: Remote procedure call protocol specification version 2, 1995.

-
- [9] Jan Stender, Björn Kolbeck, Felix Hupfeld, Eugenio Cesario, Erich Focht, Matthias Hess, Jesús Malo, and Jonathan Martí. Striping without sacrifices: maintaining posix semantics in a parallel file system. In *LASCO'08: First USENIX Workshop on Large-Scale Computing*, pages 1–8, Berkeley, CA, USA, 2008. USENIX Association.
 - [10] Osamu Tatebe, Noriyuki Soda, Youhei Morita, Satoshi Matsuoka, and Satoshi Sekiguchi. Gfarm v2: A grid file system that supports high-performance distributed and parallel data computing. In *Proceedings of the 2004 Computing in High Energy and Nuclear Physics (CHEP04)*, 2004.
 - [11] Matt Welsh, David Culler, and Eric Brewer. Seda: an architecture for well-conditioned, scalable internet services. *SIGOPS Oper. Syst. Rev.*, 35(5):230–243, 2001.