



# The XtreamFS Installation and User Guide

Version 1.1



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 \(www.XtreemFS.org\)](http://www.XtreemFS.org).

This document is © 2009 by Björn Kolbeck, Jan Stender, Minor Gordon, Felix Hupfeld, Juan Gonzales. All rights reserved.

# Contents

<b>1</b>	<b>Quick Start</b>	<b>vii</b>
<b>2</b>	<b>What is XtreamFS</b>	<b>1</b>
2.1	About XtreamFS . . . . .	1
2.2	XtreamFS Architecture . . . . .	2
2.2.1	The Components of XtreamFS . . . . .	2
2.2.2	Security . . . . .	3
2.3	Policies . . . . .	3
2.3.1	OSD and Replica Selection Policies . . . . .	3
	Attributes. . . . .	4
	Predefined Policies. . . . .	4
2.3.2	Striping Policies . . . . .	6
2.3.3	Authorization - Access Policies . . . . .	6
2.3.4	Pluggable Policies . . . . .	6
<b>3</b>	<b>XtreamFS Services</b>	<b>9</b>
3.1	Installation . . . . .	9
3.1.1	Prerequisites . . . . .	9
3.1.2	Installing from Pre-Packaged Releases . . . . .	9
3.1.3	Installing from Sources . . . . .	10
3.2	Configuration . . . . .	10
3.2.1	A Word about UUIDs . . . . .	10
3.2.2	Automatic DIR Service Discovery . . . . .	10
3.2.3	Authentication . . . . .	11
	UNIX uid/gid - NullAuthProvider . . . . .	11
	Plain SSL Certificates - SimpleX509AuthProvider . . . . .	11
	XtreamOS Certificates - XOSAuthProvider . . . . .	12
3.2.4	List of Configuration Options . . . . .	12

admin_password <i>optional</i>	12
authentication_provider	12
capability_secret	12
checksums.enabled	12
checksums.algorithm	13
database.dir	13
database.log	13
debug.level <i>optional</i>	14
debug.categories <i>optional</i>	15
dir_service.host	15
dir_service.port	15
discover <i>optional</i>	16
geographic_coordinates	16
hostname <i>optional</i>	16
http_port	16
listen.address <i>optional</i>	16
listen.port	17
local_clock_renewal	17
no_atime	17
no_fsync <i>optional</i>	17
object_dir	18
osd_check_interval	18
remote_time_sync	18
report_free_space	18
ssl.enabled	18
ssl.service_creds	19
ssl.service_creds.container	19
ssl.service_creds.pw	19
ssl.trusted_certs	19
ssl.trusted_certs.container	19
ssl.trusted_certs.pw	20
uuid	20
3.2.5 Configuring SSL Support	20
Converting PEM files to PKCS#12	20
Importing trusted certificates from PEM into a JKS	21
Sample Setup	21
3.3 Management	23

3.3.1	Starting and Stopping the XtreamFS services . . . . .	23
3.3.2	Web-based Status Page . . . . .	23
3.3.3	Creating Volumes . . . . .	23
3.3.4	Deleting Volumes . . . . .	24
<b>4</b>	<b>The XtreamFS Client</b>	<b>25</b>
4.1	Installation . . . . .	25
4.1.1	Prerequisites . . . . .	25
4.1.2	Installing from Pre-Packaged Releases . . . . .	25
4.1.3	Installing from Sources . . . . .	25
4.2	Mounting and Un-mounting . . . . .	26
<b>5</b>	<b>XtreamFS Tools</b>	<b>27</b>
5.1	Installation . . . . .	27
5.1.1	Prerequisites . . . . .	27
5.1.2	Installing from Pre-Packaged Releases . . . . .	27
5.1.3	Installing from Sources . . . . .	28
5.2	Maintenance Tools . . . . .	28
5.2.1	MRC Database Conversion . . . . .	28
5.2.2	Scrubbing and Cleanup . . . . .	29
5.3	User Tools . . . . .	29
5.3.1	Reading XtreamFS-specific File Info . . . . .	29
5.3.2	Changing Striping Policies . . . . .	30
5.3.3	Changing OSD and Replica Selection Policies . . . . .	31
5.3.4	Setting and Listing Policy Attributes . . . . .	31
5.3.5	Manual Read-Only Replication . . . . .	32
5.3.6	Automatic On-Close Replication . . . . .	33
<b>6</b>	<b>Troubleshooting and Support</b>	<b>35</b>
6.1	Logfiles . . . . .	35
6.2	Support . . . . .	35
6.3	Troubleshooting . . . . .	35
<b>A</b>	<b>XtreamOS Integration</b>	<b>37</b>
	XtreamFS Security Preparations . . . . .	37
<b>B</b>	<b>Command Line Utilities</b>	<b>39</b>



# Chapter 1

## Quick Start

This is the very short version to help you set up a local installation of XtreamFS.

1. Download XtreamFS RPMs/DEBs and install
  - (a) Download the RPMs or DEBs for your system from [the XtreamFS web-site](#)
  - (b) open a root console (su or sudo)
  - (c) install with `rpm -Uhv xtreamfs-client-1.1.x.rpm xtreamfs-server-1.1.x.rpm`
2. Start the Directory Service:  
`/etc/init.d/xtreamfs-dir start`
3. Start the Metadata Server:  
`/etc/init.d/xtreamfs-mrc start`
4. Start the OSD:  
`/etc/init.d/xtreamfs-osd start`
5. If not already loaded, load the FUSE kernel module:  
`modprobe fuse`
6. Depending on your distribution, you may have to add users to a special group to allow them to mount FUSE file systems. In openSUSE users must be in the group `trusted`, in Ubuntu in the group `fuse`. You may need to log out and log in again for the new group membership to become effective.
7. You can now close the root console and work as a regular user.
8. Wait a few seconds for the services to register at the directory service. You can check the registry by opening the DIR status page in your favorite web browser <http://localhost:30638>.
9. Create a new volume with the default settings:  
`xtfs_mkvol localhost/myVolume`
10. Create a mount point:  
`mkdir ~/xtreamfs`

11. Mount XtreamFS on your computer:

```
xtfs_mount localhost/myVolume ~/xtreamfs
```

12. Have fun ;-)

13. To un-mount XtreamFS:

```
xtfs_umount ~/xtreamfs
```

You can also mount this volume on remote computers. First make sure that the ports 32636, 32638 and 32640 are open for incoming TCP connections. You must also specify a hostname that can be resolved by the remote machine! This hostname has to be used instead of localhost when mounting.



## Chapter 2

# What is XtreamFS

### 2.1 About XtreamFS

With XtreamFS you are about to install a modern *distributed and replicated file system*. As a distributed file system, XtreamFS stores your file data on several servers and you can simply scale your file system by adding more hosts. XtreamFS is a full-featured file system that supports the full POSIX file interface, including *extended attributes* (xattrs). In case of concurrent access by several distributed programs, XtreamFS provides you currently with NFS close-to-open consistency.

Since version 1.0, XtreamFS also supports replication of files. The so called read-only replication allows you to have multiple copies of immutable files. XtreamFS also supports partial replicas which helps you to save disk capacity and network bandwidth; only data that is requested by clients is stored in partial replicas.

XtreamFS has been designed for deployment in *wide-area environments* connected by the Internet. This means that it allows you to mount an XtreamFS volume from any location, given the right permissions; but it also implies that file system installations can span multiple locations or data centers.

In a normal UNIX environment, XtreamFS has full permission and *POSIX ACL* support. XtreamFS can also be integrated into *X.509-based security architectures*. Access policies (as well several other policies) are pluggable and can be easily extended. If you deploy XtreamFS as part of an *XtreamOS* installation, you will benefit from its transparent integration with the *XtreamOS Virtual Organization (VO)* infrastructure in the form of dynamic user mappings and automatic mounting of home volumes.

If you need *high-performance* access to your files, XtreamFS can help you with support for *file striping*: XtreamFS can store a file across several storage servers and *access* the parts *in parallel*. The size of an individual stripe and the number of storage servers used can be configured on a per-file or per-directory basis.

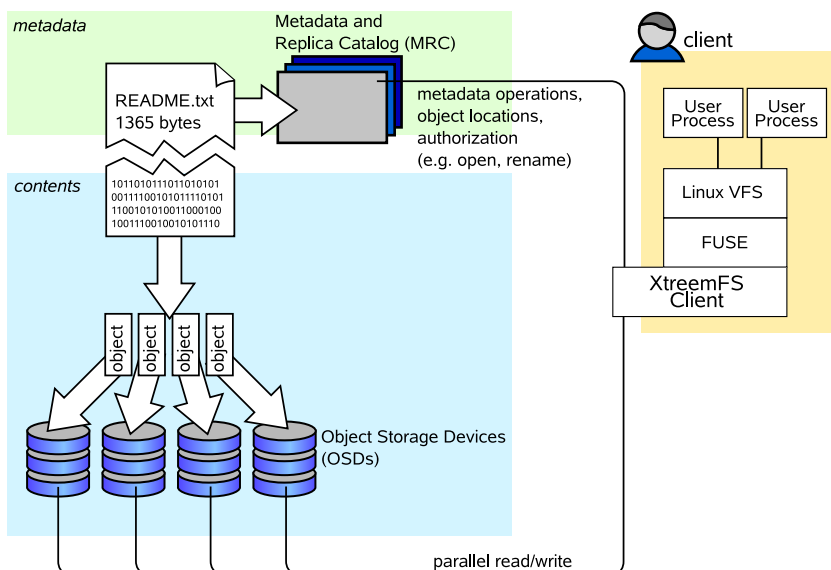


Figure 2.1: The XtremFS architecture and components.

## 2.2 XtremFS Architecture

XtremFS implements an *object-based file system architecture* (Fig. 2.1). The name of this architecture comes from the fact that an object-based file system splits file content into a series of fixed-size *objects* and stores them on its storage servers. In contrast to block-based file systems, the size of such an object can vary from file to file.

The *metadata* of a file (such as the file name or file size) is stored separate from the file content on a Metadata server. This metadata server organizes file system metadata as a set of *volumes*, each of which implements a separate file system namespace in form of a directory tree.

### 2.2.1 The Components of XtremFS

An XtremFS installation contains three types of servers that can run on one or several machines (Fig. 2.1):

- **DIR - Directory Service**  
The directory service is the central registry for all services in XtremFS. The MRC uses it to discover storage servers.
- **MRC - Metadata and Replica Catalog**  
The MRC stores the directory tree and file metadata such as file name, size or modification time. Moreover, the MRC authenticates users and authorizes access to files.
- **OSD - Object Storage Device**  
An OSD stores arbitrary objects of files; clients read and write file data on OSDs.

These servers are connected by the *client* to a file system. A client *mounts* one of the volumes of the MRC in a local directory. It translates file system calls into RPCs sent to the respective servers.

The client is implemented as a *FUSE user-level driver* that runs as a normal process. FUSE itself is a kernel-userland hybrid that connects the user-land driver to Linux' *Virtual File System (VFS)* layer where file system drivers usually live.

### 2.2.2 Security

As usual, XtremFS security differentiates between authentication and authorization. *Authentication* is the process of verifying a user's or client's identity, e.g. validating and reading an X.509 certificate. In contrast, *authorization* is the process of checking if a user has the permission to execute a certain operation, e.g. write access to a file.

By default, XtremFS uses unauthenticated and unencrypted TCP connections. However, SSL can be enabled in all XtremFS services and the client. Using SSL requires that all users and services provide valid X.509 certificates. Any data sent over a SSL connection is encrypted. Using SSL, however, will increase the resource consumption of all components, especially for connection setup (SSL handshake).

## 2.3 Policies

Many facets of the behavior of XtremFS can be configured by means of policies. A policy defines how a certain task is performed, e.g. how the MRC selects a set of OSDs for a new file, or how it distinguishes between an authorized and an unauthorized user when files are accessed. Various policies have been defined that cover different aspects.

### 2.3.1 OSD and Replica Selection Policies

When a new file is created or a replica is automatically added to a file, the MRC must decide on a set of OSDs for storing the file content. OSD Selection Policies are used to select the most suitable subset among all known OSDs.

Replica selection is a related problem. When a client opens a file with more than one replica, the MRC uses a replica selection policy to sort the list of replicas for the client. Clients will always use the first replica in the list. If a replica is not available, the client will automatically select the next replica from the list. Replica selection policies can be used to sort the replica lists, e.g. to ensure that clients use replicas which are close to them.

Both OSD and replica selection policies share a common mechanism. Input parameters of a basic policy are a set of OSDs, the current replica locations list of the file, and the IP address of the client on behalf of whom the policy was called. The output parameter is a filtered and potentially sorted subset of OSDs. Since OSD lists returned by one basic policy can be used as input parameters by another one, basic policies can be chained to define more complex composite policies.

OSD and replica selection policies are assigned at volume granularity. For further details on how to set such policies, please refer to Sec. 5.3.3.

**Attributes.** The behavior of basic policies can be further refined by means of policy attributes. Policy attributes are extended attributes with a name starting with `xtreemfs.policies.`, such as `xtreemfs.policies.minFreeCapacity`. Each time a policy attribute is set, all policies will be notified about the change. How an attribute change affects the policy behavior depends on the policy implementation.

**Predefined Policies.** The following predefined policies exist:

*Filtering policies:*

- **Default OSD filter (policy ID 1000)**  
Removes OSDs from the list that are either dead or do not have sufficient space. By default, the lower space limit for an OSD is 2GB, and the upper response time limit is 5 minutes.

Attributes:

- *free\_capacity\_bytes*: the lower space limit in bytes
- *offline\_time\_secs*: the upper response time limit in

- **FQDN-based filter (policy ID 1001)**  
Removes OSDs from the list that do not match any of the domains in a given set. By default, the set of domains contains '\*', which indicates that no domains are removed.

Attributes:

- *domains*: a comma or space-separated list of domain names. The list may include leading and trailing '\*', which will be regarded as wildcard characters.

*Grouping policies:*

- **Data center map-based grouping (policy ID 2000)**  
Removes all OSDs from the OSD set that have been used in the file's replica locations list already and selects the subset of OSDs that is closest to the client and provides enough OSDs for the new replica in a single data center.

This policy uses a statically configured datacenter map that describes the distance between datacenters. It works only with IPv4 addresses at the moment. Each datacenter has a list of matching IP addresses and networks which is used to assign clients and OSDs to datacenters. Machines in the same datacenter have a distance of 0.

This policy requires a datacenter map configuration file in `/etc/xos/xtreemfs/datacentermap` on the MRC machine which is loaded at MRC startup. This config file must contain the following parameters:

- `datacenters=A,B,C`  
A comma separated list of datacenters. Datacenter names may only contain a-z, A-Z, 0-9 and `_`.
- `distance.A-B=100`  
For each pair of datacenters, the distance must be specified. As distances are symmetric, it is sufficient to specify A to B.
- `addresses.A=192.168.1.1,192.168.2.0/24`  
For each datacenter a list of matching IP addresses or networks must be specified.
- `max_cache_size=1000`  
Sets the size of the address cache that is used to lookup IP-to-datacenter matches.

A sample datacenter map could look like this:

```
datacenters=BERLIN,LONDON,NEW_YORK
distance.BERLIN-LONDON=10
distance.BERLIN-NEW_YORK=140
distance.LONDON-NEW_YORK=110
addresses.BERLIN=192.168.1.0/24
addresses.LONDON=192.168.2.0/24
addresses.NEW_YORK=192.168.3.0/24,192.168.100.0/25
max_cache_size=100
```

- **FQDN-based grouping (policy ID 2001)**  
Removes all OSDs from the OSD set that have been used in the file's replica locations list already and selects the subset of OSDs that is closest to the client and provides enough OSDs for the new replica in a single domain.  
  
This policy uses domain names of clients and OSDs to determine the distance between a client and an OSD, as well as if OSDs are in the same domain.

*Sorting policies:*

- **Shuffling (policy ID 3000)**  
Shuffles the given list of OSDs.
- **Data center map-based sorting (policy ID 3001)**  
Sorts the list of OSDs in ascending order of their distance to the client, according to the data center map.
- **DNS based OSD Selection (policy ID 3002)**  
The FQDN of the client and all OSDs is compared and the maximum match (from the end of the FQDN) is used to sort the OSDs. The policy sorts the list of OSDs in descending order by the number of characters that match. This policy can be used to automatically select OSDs which are close to the client, if the length of the match between two DNS entries also indicate a low latency between two machines.

### 2.3.2 Striping Policies

XtreemFS allows the content of a file to be distributed among several storage devices (OSDs). This has the benefit, that the file can be read or written in parallel on multiple servers which increases the bandwidth. The more OSDs are used, the higher the bandwidth available for reading or writing. The number of OSDs is called the striping width.

A striping policy is a rule that defines how the objects are distributed on the available OSDs. Currently, XtreemFS implements only the RAID0 policy which simply stores the objects in a round robin fashion on the OSDs. The RAID0 policy has two parameters. The striping width defines to how many OSDs the file is distributed. The stripe size defines the size of each object.

Striping over several OSDs enhances the read and write bandwidth of a file, the bandwidth increases the larger the striping width. Please note, that striping also increases the probability of data loss. A striped file will become corrupted even if a single OSDs it is stored on has a disk crash.

### 2.3.3 Authorization - Access Policies

User authorization is managed by means of Access Policies. An access policy defines the access rights for any user on any file or directory contained in a volume. When creating a new volume, the access policy has to be chosen, which cannot be changed in the future. Various access policies can be used:

- Authorize All Policy (policy Id 1)  
No authorization - everyone can do everything. This policy is useful if performance of metadata operations matters more than security, since no recursive evaluation of access policies is done.
- POSIX ACLs & Permissions (policy Id 2)  
This access policy implements the traditional POSIX permissions commonly used on Linux, as well as POSIX ACLs, an extension that provides for access control at the granularity of single users and groups. POSIX permissions should be used as the default, as it guarantees maximum compatibility with other file systems.
- Volume ACLs (policy Id 3)  
Volume ACLs provide an access control model similar to POSIX ACLs & Permissions, but only allow one ACL for the whole volume. This means that there is no recursive evaluation of access rights which yields a higher performance at the price of a very coarse-grained access control.

### 2.3.4 Pluggable Policies

Administrators may extend the set of existing policies by defining *plug-in policies*. Such policies are Java classes that implement a predefined policy interface. Currently, the following policy interfaces exist:

- `org.xtreemfs.common.auth.AuthenticationProvider`  
can be used to implement an individual mechanism to authenticate users and groups
- `org.xtreemfs.mrc.ac.FileAccessPolicy`  
can be used to implement an individual access control model on files, directories and volumes
- `org.xtreemfs.mrc.osdselection.OSDSelectionPolicy`  
can be used to implement an individual policy for allocating OSDs to newly created files

Note that there may only be one authentication provider per MRC, while file access policies and OSD selection policies may differ for each volume. The former one is identified by means of its class name (property `authentication_provider`, see Sec. 3.2.4), while volume-related policies are identified by ID numbers. It is therefore necessary to add a member field

```
public static final long POLICY_ID = 4711;
```

to all such policy implementations, where 4711 represents the individual ID number. Administrators have to ensure that such ID numbers neither clash with ID numbers of built-in policies (1-9), nor with ID numbers of other plug-in policies. When creating a new volume, IDs of plug-in policies may be used just like built-in policy IDs.

Plug-in policies have to be deployed in the directory specified by the MRC configuration property `policy_dir`. The property is optional; it may be omitted if no plug-in policies are supposed to be used. An implementation of a plug-in policy can be deployed as a Java source or class file located in a directory that corresponds to the package of the class. Library dependencies may be added in the form of source, class or JAR files. JAR files have to be deployed in the top-level directory. All source files in all subdirectories are compiled at MRC start-up time and loaded on demand.





## Chapter 3

# XtreemFS Services

### 3.1 Installation

When installing XtreemFS server components, you can choose from two different installation sources: you can download one of the *pre-packaged releases* that we create for most Linux distributions or you can install directly from the *source tarball*. In the pre-packaged release, the server and the client parts are split into separate packages.

#### 3.1.1 Prerequisites

For the pre-packaged release, you will need Sun Java JRE 1.6.0 or newer to be installed on the system.

When building XtreemFS directly from the source, you need a Sun Java JDK 1.6.0 or newer, Ant 1.6.5 or newer and gmake.

#### 3.1.2 Installing from Pre-Packaged Releases

On RPM-based distributions (RedHat, Fedora, SuSE, Mandriva, XtreemOS) you can install the package with

```
$> rpm -i xtreemfs-server-1.1.x.rpm
```

For Debian-based distributions, please use the .deb package provided and install it with

```
$> dpkg -i xtreemfs-server-1.1.x.deb
```

Both packages will also install `init.d` scripts for an automatic start-up of the services. Use `insserv xtreemfs-dir`, `insserv xtreemfs-mrc` and `insserv xtreemfs-osd`, respectively, to automatically start the services during boot.

### 3.1.3 Installing from Sources

Extract the tarball with the sources. Change to the top level directory and execute

```
$> make server
```

This will build the XtreamFS server and Java-based tools. When done, execute

```
$> sudo make install
```

to complete the installation.

## 3.2 Configuration

Generally, the configuration files of XtreamFS are located in `/etc/xos/xtreemfs/` if you installed from packages.

### 3.2.1 A Word about UUIDs

XtreamFS uses UUIDs (Universally Unique Identifiers) to be able to identify services and their associated state independently from the machine they are installed on. This implies that you cannot change the uuid of a MRC or OSD after it has been used for the first time!

The Directory Service keeps a mapping from UUID to a port number and IP address or hostname. Currently, each UUID can only be assigned to a single endpoint; the netmask must be “\*” which means that this mapping is valid in all networks. Upon first start-up, OSDs and MRCs will create the mapping if it does not exist. They will use the first available network device with a public address.

Changing the IP address, hostname or port is possible at any time. Due to the caching of UUIDs in all components it can take some time until the new UUID mapping is used by all OSDs, MRCs and clients. The TTL defines how long an XtreamFS component is allowed to keep entries cached. The default value is 3600 seconds (1 hour). It should be set to shorter durations if services change their IP address frequently.

To create a globally unique UUID you can use tools like `uuidgen`. During installation the post-install script will automatically create a UUID for each OSD and MRC if it does not have a UUID assigned.

### 3.2.2 Automatic DIR Service Discovery

The OSD and MRC discover the DIR service automatically by default. On startup they will broadcast requests to the local LAN and wait up to 10s for a response from a DIR. This works only in a local LAN environment as broadcast messages are not routed to other networks. Local firewalls on the computers on which the services are running can also prevent the automatic discovery from working. The services will

select the first DIR which responded which can lead to non-deterministic behavior if multiple DIR services are present.

**Security:** The automatic discovery is a potential security risk when used in untrusted environments as any user can start-up DIR services.

A statically configured DIR address and port can be used to disable DIR discovery in the OSD and MRC (see Sec. 3.2.4, `dir_service`). By default, the DIR responds to UDP broadcasts. To disable this feature, set `discover = false` in the DIR service config file.

### 3.2.3 Authentication

XtreemFS has an interface which allows MRC administrators to choose the way of authenticating users. Basically, an MRC has two sources of information on users. The first one is the user id and group ids sent by the client along with each request. In addition, the MRC can use information included in the certificates if SSL is enabled. The Authentication Providers are modules that implement different methods for retrieving the user and group IDs to use.

#### UNIX uid/gid - NullAuthProvider

The NullAuthProvider is the default Authentication Provider. It simply uses the user ID and group IDs sent by the XtreemFS client. This means that the client is trusted to send the correct user/group IDs.

The XtreemFS Client will send the user ID and group IDs of the process which executed the file system operation, not of the user who mounted the volume!

The superuser is identified by the user ID `root` and is allowed to do everything on the MRC. This behavior is similar to NFS with `no_root_squash`.

#### Plain SSL Certificates - SimpleX509AuthProvider

XtreemFS supports two X.509 certificate “types” which can be used by the client. When mounted with a service/host certificate the XtreemFS client is regarded as a trusted system component. The MRC will accept any user ID and groups sent by the client and use them for authorization as with the NullAuthProvider. This setup is useful for volumes which are used by multiple users.

The second certificate type are regular user certificates. The MRC will only accept the user name and group from the certificate and ignore the user ID and groups sent by the client. Such a setup is useful if users are allowed to mount XtreemFS from untrusted machines.

Both certificates are regular X.509 certificates. Service and host certificates are identified by a Common Name (CN) starting with `host/` or `xtreemfs-service/`, which can easily be used in existing security infrastructures. All other certificates are assumed to be user certificates.

If a user certificate is used, XtreemFS will take the Distinguished Name (DN) as the user ID and the Organizational Unit (OU) as the group ID.

Superusers must have `xtreemfs-admin` as part of their Organizational Unit (OU).

### XtreemOS Certificates - XOSAuthProvider

In contrast to plain X.509 certificates, XtreemOS embeds additional user information as extensions in XtreemOS-User-Certificates. This authentication provider uses this information (global UID and global GIDs), but the behavior is similar to the SimpleX509AuthProvider.

The superuser is identified by being member of the VOAdmin group.

#### 3.2.4 List of Configuration Options

All configuration parameters that may be used to define the behavior of the different services are listed in this section. Unless marked as optional, a parameter has to occur (exactly once) in a configuration file.

##### `admin_password` *optional*

Services	DIR, MRC, OSD
Values	String
Default	empty
Description	Defines the admin password that must be sent to authorize requests like volume creation, deletion or shutdown.

##### `authentication_provider`

Services	MRC
Values	Java class name
Default	<code>org.xtreemfs.common.auth.NullAuthProvider</code>
Description	Defines the Authentication Provider to use to retrieve the user identity (user ID and group IDs). See Sec. 3.2.3 for details.

##### `capability_secret`

Services	MRC, OSD
Values	String
Default	-
Description	Defines a shared secret between the MRC and all OSDs. The secret is used by the MRC to sign capabilities, i.e. security tokens for data access at OSDs. In turn, an OSD uses the secret to verify that the capability has been issued by the MRC.

##### `checksums.enabled`

Services	OSD
Values	true, false
Default	false
Description	If set to true, the OSD will calculate and store checksums for newly created objects. Each time a checksummed object is read, the checksum will be verified.

**checksums.algorithm**

Services	OSD
Values	Adler32, CRC32
Default	Adler32
Description	Must be specified if <code>checksums.enabled</code> is enabled. This property defines the algorithm used to create OSD checksums.

**database.dir**

Services	DIR, MRC
Values	absolute file system path to a directory
Default	DIR: <code>/var/lib/xtreemfs/dir/database</code> , MRC: <code>/var/lib/xtreemfs/mrc/database</code>
Description	The directory in which the Directory Service or MRC will store their databases. This directory should never be on the same partition as any OSD data, if both services reside on the same machine. Otherwise, deadlocks may occur if the partition runs out of free disk space!

**database.log**

Services	MRC
Values	absolute file system path
Default	MRC: <code>/var/lib/xtreemfs/mrc/dblog</code>
Description	The directory the MRC uses to store database logs. This directory should never be on the same partition as any OSD data, if both services reside on the same machine. Otherwise, deadlocks may occur if the partition runs out of free disk space!

`debug.level` *optional*

Services	DIR, MRC, OSD
Values	0, 1, 2, 3, 4, 5, 6, 7
Default	6
Description	<p>The debug level determines the amount and detail of information written to logfiles. Any debug level includes log messages from lower debug levels. The following log levels exist:</p> <ul style="list-style-type: none"><li>0 - fatal errors</li><li>1 - alert messages</li><li>2 - critical errors</li><li>3 - normal errors</li><li>4 - warnings</li><li>5 - notices</li><li>6 - info messages</li><li>7 - debug messages</li></ul>

**debug.categories** *optional*

Services	DIR, MRC, OSD
Values	all, lifecycle, net, auth, stage, proc, db, misc
Default	all
Description	Debug categories determine the domains for which log messages will be printed. By default, there are no domain restrictions, i.e. log messages from all domains will be included in the log. The following categories can be selected:
	all - no restrictions on the category
	lifecycle - service lifecycle-related messages, including startup and shut-down events
	net - messages pertaining to network traffic and communication between services
	auth - authentication and authorization-related messages
	stage - messages pertaining to the flow of requests through the different stages of a service
	proc - messages about the processing of requests
	db - messages that are logged in connection with database accesses
	misc - any other log messages that do not fit in one of the previous categories
	Note that it is possible to specify multiple categories by means of a comma or space-separated list.

**dir\_service.host**

Services	MRC, OSD
Values	hostname or IP address
Default	localhost
Description	Specifies the hostname or IP address of the directory service (DIR) at which the MRC or OSD should register. The MRC also uses this directory service to find OSDs. If set to <code>.autodiscover</code> the service will use the automatic DIR discovery mechanism (see Sec. 3.2.2).

**dir\_service.port**

Services	MRC, OSD
Values	1 .. 65535
Default	32638
Description	Specifies the port on which the remote directory service is listening. Must be identical to the <code>listen_port</code> in your directory service configuration.

**discover** *optional*

Services	DIR
Values	true, false
Default	true
Description	If set to true the DIR will received UDP broadcasts and advertise itself in response to XtremFS components using the DIR automatic discovery mechanism. If set to false, the DIR will ignore all UDP traffic. For details see Sec. 3.2.2.

**geographic\_coordinates**

Services	DIR, MRC, OSD
Values	String
Default	empty
Description	Specifies the geographic coordinates which are registered with the directory service. Used e.g. by the web console.

**hostname** *optional*

Services	MRC, OSD
Values	String
Default	-
Description	If specified, it defines the host name that is used to register the service at the directory service. If not specified, the host address defined in <code>listen.address</code> will be used if specified. If neither <code>hostname</code> nor <code>listen.address</code> are specified, the service itself will search for externally reachable network interfaces and advertise their addresses.

**http\_port**

Services	DIR, MRC, OSD
Values	1 .. 65535
Default	30636 (MRC), 30638 (DIR), 30640 (OSD)
Description	Specifies the geographic coordinates which are registered with the directory service. Used e.g. by the web console.

**listen.address** *optional*

Services	OSD
Values	IP address
Default	-
Description	If specified, it defines the interface to listen on. If not specified, the service will listen on all interfaces (any).



**listen.port**

Services	DIR, MRC, OSD
Values	1 .. 65535
Default	DIR: 32638, MRC: 32636, OSD: 32640
Description	The port to listen on for incoming connections (TCP). The OSD uses TCP and UDP on the specified port. Make sure to configure your firewall to allow incoming TCP and UDP traffic on the specified port.

**local\_clock\_renewal**

Services	MRC, OSD
Values	milliseconds
Default	50
Description	Reading the system clock is a slow operation on some systems (e.g. Linux) as it is a system call. To increase performance, XtremFS services use a local variable which is only updated every <code>local_clock_renewal</code> milliseconds.

**no\_atime**

Services	MRC
Values	true, false
Default	true
Description	The POSIX standard defines that the atime (timestamp of last file access) is updated each time a file is opened, even for read. This means that there is a write to the database and hard disk on the MRC each time a file is read. To reduce the load, many file systems (e.g. ext3) including XtremFS can be configured to skip those updates for performance. It is strongly suggested to disable atime updates by setting this parameter to true.

**no\_fsync *optional***

Services	MRC
Values	true, false
Default	false
Description	By default, the MRC will write all file-modifying operations (such as create file, delete etc.) to disk followed by a <code>fsync</code> to ensure data is written to the hard disk. While this ensures maximum data safety in case of crash of the MRC server, it also reduces the performance of the MRC. Set this to true, if you want much higher performance at the risk of losing some recent file operations in case of a server crash.

**object\_dir**

Services	OSD
Values	absolute file system path to a directory
Default	/var/lib/xtreemfs/osd/
Description	The directory in which the OSD stores the objects. This directory should never be on the same partition as any DIR or MRC database, if both services reside on the same machine. Otherwise, deadlocks may occur if the partition runs out of free disk space!

**osd\_check\_interval**

Services	MRC
Values	seconds
Default	300
Description	The MRC regularly asks the directory service for suitable OSDs to store files on (see OSD Selection Policy, Sec. 2.3.1). This parameter defines the interval between two updates of the list of suitable OSDs.

**remote\_time\_sync**

Services	MRC, OSD
Values	milliseconds
Default	30,000
Description	MRCs and OSDs all synchronize their clocks with the directory service to ensure a loose clock synchronization of all services. This is required for leases to work correctly. This parameter defines the interval in milliseconds between time updates from the directory service.

**report\_free\_space**

Services	OSD
Values	true, false
Default	true
Description	If set to true, the OSD will report its free space to the directory service. Otherwise, it will report zero, which will cause the OSD not to be used by the OSD Selection Policies (see Sec. 2.3.1).

**ssl.enabled**

Services	DIR, MRC, OSD
Values	true, false
Default	false
Description	If set to true, the service will use SSL to authenticate and encrypt connections. The service will not accept non-SSL connections if <code>ssl.enabled</code> is set to true.

**ssl.service\_creds**

Services	DIR, MRC, OSD
Values	path to file
Default	DIR: /etc/xos/xtreemfs/truststore/certs/ds.p12, MRC: /etc/xos/xtreemfs/truststore/certs/mrc.p12, OSD: /etc/xos/xtreemfs/truststore/certs/osd.p12
Description	Must be specified if <code>ssl.enabled</code> is enabled. Specifies the file containing the service credentials (X.509 certificate and private key). PKCS#12 and JKS format can be used, set <code>ssl.service_creds.container</code> accordingly. This file is used during the SSL handshake to authenticate the service.

**ssl.service\_creds.container**

Services	DIR, MRC, OSD
Values	pkcs12 or JKS
Default	pkcs12
Description	Must be specified if <code>ssl.enabled</code> is enabled. Specifies the file format of the <code>ssl.service_creds</code> file.

**ssl.service\_creds.pw**

Services	DIR, MRC, OSD
Values	String
Default	-
Description	Must be specified if <code>ssl.enabled</code> is enabled. Specifies the password which protects the credentials file <code>ssl.service_creds</code> .

**ssl.trusted\_certs**

Services	DIR, MRC, OSD
Values	path to file
Default	/etc/xos/xtreemfs/truststore/certs/xosrootca.jks
Description	Must be specified if <code>ssl.enabled</code> is enabled. Specifies the file containing the trusted root certificates (e.g. CA certificates) used to authenticate clients.

**ssl.trusted\_certs.container**

Services	DIR, MRC, OSD
Values	pkcs12 or JKS
Default	JKS
Description	Must be specified if <code>ssl.enabled</code> is enabled. Specifies the file format of the <code>ssl.trusted_certs</code> file.

**ssl.trusted\_certs.pw**

Services	DIR, MRC, OSD
Values	String
Default	-
Description	Must be specified if <code>ssl.enabled</code> is enabled. Specifies the password which protects the trusted certificates file <code>ssl.trusted_certs</code> .

**uuid**

Services	MRC, OSD
Values	String, but limited to alphanumeric characters, - and .
Default	-
Description	Must be set to a unique identifier, preferably a UUID according to RFC 4122. UUIDs can be generated with <code>uuidgen</code> . Example: <code>eacb6bab-f444-4ebf-a06a-3f72d7465e40</code> .

**3.2.5 Configuring SSL Support**

In order to enable certificate-based authentication in an XtremFS installation, services need to be equipped with X.509 certificates. Certificates are used to establish a mutual trust relationship among XtremFS services and between the XtremFS client and XtremFS services.

It is not possible to mix SSL-enabled and non-SSL services in an XtremFS installation!

Each XtremFS service needs a certificate and a private key in order to be run. Once they have been created and signed, the credentials may need to be converted into the correct file format. XtremFS services also need a *trust store* that contains all trusted Certification Authority certificates.

By default, certificates and credentials for XtremFS services are stored in

```
/etc/xos/xtreemfs/truststore/certs
```

**Converting PEM files to PKCS#12**

The simplest way to provide the credentials to the services is by converting your signed certificate and private key into a PKCS#12 file using `openssl`:

```
$> openssl pkcs12 -export -in ds.pem -inkey ds.key \
    -out ds.p12 -name "DS"
$> openssl pkcs12 -export -in mrc.pem -inkey mrc.key \
    -out mrc.p12 -name "MRC"
$> openssl pkcs12 -export -in osd.pem -inkey osd.key \
    -out osd.p12 -name "OSD"
```

This will create three PKCS12 files (`ds.p12`, `mrc.p12` and `osd.p12`), each containing the private key and certificate for the respective service. The passwords chosen when asked must be set as a property in the corresponding service configuration file.

### Importing trusted certificates from PEM into a JKS

The certificate (or multiple certificates) from your CA (or CAs) can be imported into a Java Keystore (JKS) using the Java keytool which comes with the Java JDK or JRE.

Execute the following steps for each CA certificate using the same keystore file.

```
$> keytool -import -alias rootca -keystore trusted.jks \
    -trustcacerts -file ca-cert.pem
```

This will create a new Java Keystore `trusted.jks` with the CA certificate in the current working directory. The password chosen when asked must be set as a property in the service configuration files.

Note: If you get the following error

```
$> keytool error: java.lang.Exception: Input not an X.509 certificate
```

you should remove any text from the beginning of the certificate (until the `---BEGIN CERTIFICATE---` line).

### Sample Setup

Users can easily set up their own CA (certificate authority) and create and sign certificates using `openssl` for a test setup.

1. Set up your test CA.

- (a) Create a directory for your CA files

```
$> mkdir ca
```

- (b) Create a private key and certificate request for your CA.

```
$> openssl req -new -newkey rsa:1024 -nodes -out ca/ca.csr \
    -keyout ca/ca.key
```

Enter something like `XtreemFS-DEMO-CA` as the common name (or something else, but make sure the name is different from the server and client name!).

- (c) Create a self-signed certificate for your CA which is valid for one year.

```
$> openssl x509 -trustout -signkey ca/ca.key -days 365 -req \
    -in ca/ca.csr -out ca/ca.pem
```

- (d) Create a file with the CA's serial number

```
$> echo "02" > ca/ca.srl
```

2. Set up the certificates for the services and the XtreemFS Client.

Replace *service* with *dir*, *mrc*, *osd* and *client*.

- (a) Create a private key for the service.

Use `XtreemFS-DEMO-service` as the common name for the certificate.

```
$> openssl req -new -newkey rsa:1024 -nodes
      -out service.req
      -keyout service.key
```

- (b) Sign the certificate with your demo CA.  
The certificate is valid for one year.

```
$> openssl x509 -CA ca/ca.pem -CAkey ca/ca.key
      -CAserial ca/ca.srl -req
      -in service.req
      -out service.pem -days 365
```

- (c) Export the service credentials (certificate and private key) as a PKCS#12 file.

Use “passphrase” as export password. You can leave the export password empty for the XtremFS Client to avoid being asked for the password on mount.

```
$> openssl pkcs12 -export -in service.pem -inkey service.key
      -out service.p12 -name "service"
```

- (d) Copy the PKCS#12 file to the certificates directory.

```
$> mkdir -p /etc/xos/xtreemfs/truststore/certs
$> cp service.p12 /etc/xos/xtreemfs/truststore/certs
```

3. Export your CA’s certificate to the trust store and copy it to the certificate dir.  
You should answer “yes” when asked “Trust this certificate”.  
Use “passphrase” as passphrase for the keystore.

```
$> keytool -import -alias ca -keystore trusted.jks \
      -trustcacerts -file ca/ca.pem
$> cp trusted.jks /etc/xos/xtreemfs/truststore/certs
```

4. Configure the services. Edit the configuration file for all your services. Set the following configuration options (see Sec. 3.2 for details).
- ```
ssl.enabled = true
ssl.service_creds.pw = passphrase
ssl.service_creds.container = pkcs12
ssl.service_creds = /etc/xos/xtreemfs/truststore/certs/service.p12
ssl.trusted_certs = /etc/xos/xtreemfs/truststore/certs/trusted.jks
ssl.trusted_certs.pw = passphrase
ssl.trusted_certs.container = jks
```

5. Start up the XtremFS services (see Sec. 3.3.1).

6. Create a new volume (see Sec. 3.3.3 for details).

```
$> xtfs_mkvol --pkcs12-file-path=\
      /etc/xos/xtreemfs/truststore/certs/client.p12 localhost/test
```

7. Mount the volume (see Sec. 4.2 for details).

```
$> xtfs_mount --pkcs12-file-path=\
      /etc/xos/xtreemfs/truststore/certs/client.p12 localhost/test /mnt
```

### 3.3 Management

This section covers all tools and functionality for XtreamFS management and tracing. In general, the use of management tools is restricted to superusers.

#### 3.3.1 Starting and Stopping the XtreamFS services

If you installed a *pre-packaged release* you can start, stop and restart the services with the `init.d` scripts:

```
$> /etc/init.d/xtreemfs-ds start
$> /etc/init.d/xtreemfs-mrc start
$> /etc/init.d/xtreemfs-osd start
```

or

```
$> /etc/init.d/xtreemfs-ds stop
$> /etc/init.d/xtreemfs-mrc stop
$> /etc/init.d/xtreemfs-osd stop
```

**Note** that the Directory Service should be started first, in order to allow other services an immediate registration. Once a Directory Service and at least one OSD and MRC are running, XtreamFS is operational.

#### 3.3.2 Web-based Status Page

The XtreamFS services all have a HTML status page which can be used to check if the service is working correctly (Fig. 3.1). It can be displayed by opening the service URL in your favorite web browser, e.g.

`http://my-mrc-host.com:30636/`. Make sure to use the right port, see `http_port` in the service config file.


#### 3.3.3 Creating Volumes

Volumes can be created with the `xtfs_mkvol` command line utility. Please see `man xtfs_mkvol` for a full list of options and usage.

When creating a volume, it is recommended to specify the access policy (see Sec. 2.3.3). If not specified, POSIX permissions/ACLs will be chosen by default. Access policies cannot be changed afterwards.

An OSD selection policy (see Sec. 2.3.1) can also be specified per volume, but can be changed anytime. By default, a random selection of available OSDs is assigned to newly created files.

In addition, it is recommended to set a default striping policy (see Sec. 2.3.2). If no per-file or per-directory default striping policy overrides the volume's default striping policy, the volume's policy is used for new files and directories. If no volume policy is explicitly defined, a RAID0 policy with a stripe size of 128kB and a width of 1 will be assigned to the volume.



OSD [osd.xtreemfs.org:32640](http://osd.xtreemfs.org:32640)

|                                  |                                                             |
|----------------------------------|-------------------------------------------------------------|
| <b>Configuration</b>             |                                                             |
| TCP & UDP port                   | 32641                                                       |
| Directory Service                | <a href="http://localhost:32638">http://localhost:32638</a> |
| Debug Level                      | 1                                                           |
| <b>Load</b>                      |                                                             |
| # HTTP connections (pinky)       | 1                                                           |
| HTTP server (pinky) queue length | 1                                                           |
| Storage Stage queue length       | 0                                                           |
| Open files                       | 0                                                           |
| <b>Transfer</b>                  |                                                             |
| # object written                 | 0                                                           |
| # object read                    | 0                                                           |
| bytes sent                       | 0 bytes                                                     |
| bytes received                   | 0 bytes                                                     |
| # GMAX packets received          | 0                                                           |
| # GMAX requests sent             | 0                                                           |
| # files deleted                  | 0                                                           |
| <b>VM Info / Memory</b>          |                                                             |
| Free Disk Space                  | 21.57 GB                                                    |
| Memory free/maxtotal             | 28.79 MB / 489.19 MB / 31.69 MB                             |
| Buffer Pool stats                | 8192: poolSize = 2 numRequests = 3 creates = 2              |
|                                  | 65536: poolSize = 0 numRequests = 0 creates = 0             |
|                                  | 524288: poolSize = 0 numRequests = 1 creates = 1            |
|                                  | 2097152: poolSize = 0 numRequests = 1 creates = 1           |
| <b>Time</b>                      |                                                             |
| global XtreemFS time             | Wed Jul 16 14:23:30 CEST 2008 (1216211010523)               |
| resync interval for global time  | 60000 ms                                                    |
| local system time                | Wed Jul 16 14:23:30 CEST 2008 (1216211010521)               |
| local time update interval       | 50 ms                                                       |

Figure 3.1: OSD status web page

An example call to `xtfs_mkvol` for creating a volume with POSIX ACLs, 256kB stripe size and a stripe width of 1 (which means no striping):

```
$> xtfs_mkvol -a POSIX -p RAID0 -s 256 -w 1 \
    my-mrc-host.com:32636/myVolume
```

### 3.3.4 Deleting Volumes

The `xtfs_rmvol` tool can be used to delete a volume. This also *deletes all files and data on that volume!* Please see `man xtfs_rmvol` for a full list of options and usage.

Example call to `xtfs_rmvol` to delete `myVolume`:

```
$> xtfs_rmvol my-mrc-host.com:32636/myVolume
```

Volume deletion is restricted to volume owners and privileged users.



## Chapter 4

# The XtreamFS Client

### 4.1 Installation

As for the XtreamFS Services, there are two different installation sources for the XtreamFS Client: *pre-packaged releases* and *source tarballs*.

#### 4.1.1 Prerequisites

For both installations you need FUSE 2.6 or newer, openssl 0.9.8 or newer and a Linux 2.6 kernel. For optimal performance we suggest to use FUSE 2.8 with a kernel version 2.6.26 or newer.

To install the client tools,

To build the XtreamFS Client from sources, you need the openssl headers (e.g. openssl-devel package), python  $\geq 2.4$ , and gcc-c++  $\geq 4.2$ .

#### 4.1.2 Installing from Pre-Packaged Releases

On RPM-based distributions (RedHat, Fedora, SuSE, Mandriva, XtreamOS) you can install the package with

```
$> rpm -i xtreamfs-client-1.1.x.rpm
```

For Debian-based distributions, please use the .deb package provided and install it with

```
$> dpkg -i xtreamfs-client-1.1.x.deb
```

#### 4.1.3 Installing from Sources

Extract the tarball with the sources. Change to the top level directory and execute

```
$> make client
```

This will build the XtremFS client and non-Java-based tools. When done, execute

```
$> sudo make install
```

to complete the installation.

## 4.2 Mounting and Un-mounting

Before mounting XtremFS volumes, please ensure that the FUSE kernel module is loaded. Please check your distribution's manual to see, if users must be in a special group (e.g. `trusted` in `openSUSE`) to be allowed to mount FUSE.

```
$> su
Password:
#> modprobe fuse
#> exit
```

To mount an XtremFS volume use the `xtfs_mount` tool.

```
$> xtfs_mount remote.dir.machine/myVolume /xtreemfs
```

`remote.dir.machine` describes the host with the Directory Service at which the volume is registered; `myVolume` is the name of the volume name to be mounted. `/xtreemfs` is the directory on the local file system to which the XtremFS volume will be mounted. For more options, please refer to `man xtfs_mount`.

The client will immediately go into background and won't display any error messages. Use the `-f` option to prevent the mount process from going into background and get all error messages printed to the console.

Access to a FUSE mount is usually restricted to the user who mounted the volume. To allow the root user or any other user on the system to access the mounted volume, the FUSE options `-o allow_root` and `-o allow_other` can be used with `xtfs_mount`. They are, however, mutually exclusive. In order to use these options, the system administrator must create a FUSE configuration file `/etc/fuse.conf` and add a line `user_allow_other`.

To check that a volume is mounted use the `mount` command. It outputs a list of all mounts in the system. XtremFS volumes are listed as type `fuse`:

```
/dev/fuse on /xtreemfs type fuse (rw,nosuid,nodev,user=userA)
```

Volumes are unmounted using the `xtfs_umount` tool.

```
$> xtfs_umount /xtreemfs
```

## Chapter 5

# XtreemFS Tools

### 5.1 Installation

When installing XtreemFS tools, you can choose from two different installation sources: you can download one of the *pre-packaged releases* that we create for most Linux distributions or you can install directly from the *source tarball*. In the pre-packaged release, the server and the client parts are split into separate packages.

#### 5.1.1 Prerequisites

For the pre-packaged release, you will need Sun Java JRE 1.6.0 or newer to be installed on the system.

When building XtreemFS directly from the source, you need a Sun Java JDK 1.6.0 or newer, Ant 1.6.5 or newer and gmake.

#### 5.1.2 Installing from Pre-Packaged Releases

On RPM-based distributions (RedHat, Fedora, SuSE, Mandriva, XtreemOS) you can install the package with

```
$> rpm -i xtreemfs-tools-1.1.x.rpm
```

For Debian-based distributions, please use the .deb package provided and install it with

```
$> dpkg -i xtreemfs-tools-1.1.x.deb
```

All XtreemFS tools will be installed to /usr/bin.

### 5.1.3 Installing from Sources

Extract the tarball with the sources. Change to the top level directory and execute

```
$> make
```

When done, execute

```
$> sudo make install
```

to complete the installation. Note that this will also install the XtreamFS client and servers.

## 5.2 Maintenance Tools

Tools for the maintenance of an XtreamFS installation will be described in the following.

### 5.2.1 MRC Database Conversion

The format in which the MRC stores its data on disk may change with future XtreamFS versions. In order that XtreamFS server components may be updated without losing the whole content of the file system, it is possible to create a version-independent XML representation of the metadata stored in MRC database.

Such an XML representation can e.g. be created as follows:

```
$> xtfs_mrcdbtool -mrc oncrpc://my-mrc-host.com:32636 \  
    dump /tmp/dump.xml
```

This call will create a file `dump.xml` containing the entire MRC database content in the `/tmp` directory at `my-mrc-host.com`.

To restore an MRC database from a dump, execute

```
$> xtfs_mrcdbtool -mrc oncrpc://my-mrc-host.com:32636 \  
    restore /tmp/dump.xml
```

This will restore the database stored in `/tmp/dump.xml` at `my-mrc-host.com`. Note that for safety reasons, it is only possible to restore a database from a dump if the database of the running MRC does not have any content. To restore an MRC database, it is thus necessary to delete all MRC database files before starting the MRC.

### 5.2.2 Scrubbing and Cleanup

In real-world environments, errors occur in the course of creating, modifying or deleting files. This can cause corruptions of file data or metadata. Such things happen e.g. if the client is suddenly terminated, or loses connection with a server component. There are several such scenarios: if a client writes to a file but does not report file sizes received from the OSD back to the MRC, inconsistencies between the file size stored in the MRC and the actual size of all objects in the OSD will occur. If a client deletes a file from the directory tree, but cannot reach the OSD, orphaned objects will remain on the OSD. If an OSD is terminated during an ongoing write operation, file content will become corrupted.

In order to detect and, if possible, resolve such inconsistencies, tools for scrubbing and OSD cleanup exist. To check the consistency of file sizes and checksums, the following command can be executed:

```
$> xtfs_scrub -dir oncrpc://my-dir-host.com:32638 xtreamfsVolume
```

This will scrub each file in the volume myVolume, i.e. check file size consistency and set the correct file size on the MRC, if necessary, and check whether an invalid checksum in the OSD indicates a corrupted file content. The `-dir` argument specifies the directory service that will be used to resolve service UUIDs. Please see `man xtfs_scrub` for further details.

A second tool scans an OSD for orphaned objects, which can be used as follows:

```
$> xtfs_cleanup -dir oncrpc://localhost:32638 \
    uuid:u2i3-28isu2-iwuv29-isjd83
```

The given UUID identifies the OSD to clean and will be resolved by the directory service defined by the `-dir` option (localhost:32638 in this example). The process will be started and can be stopped by setting the option `-stop`. To watch the cleanup progress use option `-i` for the interactive mode. For further information see `man xtfs_cleanup`.

## 5.3 User Tools

There is a range of tools for the specific features of XtreamFS, which will be described in the following.

### 5.3.1 Reading XtreamFS-specific File Info

In addition to the regular file system information provided by the `stat` Linux utility, XtreamFS provides the `xtfs_stat` tool which displays XtreamFS specific information for a file or directory.

```
$> cd /xtreamfs
$> echo 'Hello World' > test.txt
$> xtfs_stat test.txt
```

will produce output similar to the following:

```
filename                test.txt
XtreemFS URI            oncrpc://localhost/test/test.txt
XtreemFS fileID         41e9a04d-0b8b-467b-94ef-74ade02a2dc9:6
object type             regular file
owner                   stender
group                   users
read-only               false

XtreemFS replica list
  list version          0
  replica update policy
  -----
  replica 1 SP          STRIPING_POLICY_RAID0, 128kb, 1
  replica 1 OSDs        [{address=127.0.0.1:32640, uuid=OSD1}]
  replica 1 repl. flags 0x1
  -----
```

The fileID is the unique identifier of the file used on the OSDs to identify the file's objects. The owner/group fields are shown as reported by the MRC, you may see other names on your local system if there is no mapping (i.e. the file owner does not exist as a user on your local machine). Finally, the XtreemFS replica list shows the striping policy of the file, the number of replicas and for each replica, the OSDs used to store the objects.

### 5.3.2 Changing Striping Policies

It is not (yet) possible to change the striping policy of an existing file, as this would require moving and reformatting data among OSDs. However, individual striping policies can be assigned to new files (i.e. empty files) by changing the default striping policy of the parent directory or volume. For this purpose, XtreemFS provides the `xtfs_sp` tool. The tool can be used to change the striping policy that will be assigned to newly created files.

```
$> xtfs_sp --set -p RAID0 -w 4 -s 256 /xtreemfs
```

In addition, the tool can be used to retrieve the default striping policy of a volume or directory.

```
$> xtfs_sp --get /xtreemfs
```

The output will be similar to the following:

```
file:           /xtreemfs
policy:         STRIPING_POLICY_RAID0
stripe-size:    4
width (kB):     256
```

When creating a new file, XtreamFS will first check whether a default striping policy has been assigned to the parent directory. If this is not the case, the default striping policy for the volume will be used as the striping policy for the new file. Changing a volume's or directory's default striping policy requires superuser access rights or ownership of the directory or volume.

### 5.3.3 Changing OSD and Replica Selection Policies

The OSD and replica selection policy can only be set for the entire volume. The policies are described in Sec. 2.3.1. To show and modify these policies for a volume, use the `xtfs_repl` tool.

```
$> xtfs_repl --rsp_get /xtreemfs
```

displays the current replica selection policy used for the volume. To change the policy to use the datacenter map for a volume, use

```
$> xtfs_repl --rsp_set dcmap /xtreemfs
```

Similarly, OSD selection policies are set and retrieved:

```
$> xtfs_repl --osp_get /xtreemfs
```

retrieves the OSD selection policy, whereas it can be set via

```
$> xtfs_repl --osp_set dcmap /xtreemfs
```

If you want to use a custom (i.e. plug-in) policy, pass the ID instead of the name.

### 5.3.4 Setting and Listing Policy Attributes

OSD and replica selection policy behavior can be defined by means of policy attributes. Policy attributes can be set as follows:

```
$> xtfs_repl --pol_attr_set domains "*.xtreemfs.org xtreemfs.de" /xtreemfs
```

A list of all policy attributes that have been set can be shown as follows:

```
$> xtfs_repl --pol_attrs_get /xtreemfs
```

### 5.3.5 Manual Read-Only Replication

Read-only replication allows you to create copies of files on several OSDs. Since the files are immutable (read-only) there is no overhead for coordinating the replicas. XtremFS supports partial and full replicas. A partial replica will only fetch the data requested by a client (on demand). A full replica does the same, but it also fetches the data in the background without client requests until all data of the file has been replicated (background replication).

With XtremFS all replicas are initially empty and can be used immediately by applications. If a replica does not have the requested data, it fetches it from another replica and saves it locally for future requests. This helps you to reduce start-up times, and saves network bandwidth and disk usage.

The XtremFS client gets a list of replicas from the MRC when opening a file. This list is sorted according to the volume's replica selection policy, i.e. the first replica in the list is the "best" replica for the client. The client will use the first replica in the list and automatically switches to the next one if it is not reachable.

So far XtremFS only supports to manage replicas manually. Before the replication can be used the file must be marked as read-only with the following command:

```
$> xtfs_repl --set_readonly local-path-of-file
```

After a file is marked as read-only replicas can be added. The tool supports different replica creation modes. The automatic mode retrieves a list of OSDs from the MRC and chooses the best OSD according to the current replica selection policy. You can also select an OSD by specifying its UUID on the command line.

By default, partial replicas will be created. To create a full replica the option `-full` must be set.

XtremFS supports different transfer strategies which has an big impact on the speed of the replication and the order in which objects are fetched. A transfer strategy must be chosen for each replica.

To create a full replica and usage of the random transfer strategy, the following command must be used. The OSDs are selected according to the volume's replica selection policy. See Sec. 2.3.1 for a description of the policies.

```
$> xtfs_repl --add_auto --full --strategy random local-path-of-file
```

To list all replicas and OSDs of the file use:

```
$> xtfs_repl -l local-path-of-file
```

Removing also supports multiple modes. Again there is a mode which is choosing the replica randomly. To remove a specific replica, the head-OSD (first OSD) of this replica must be given as an argument. At least one complete replica must exist, so the tool will not remove a complete replica, if no others exists. To remove the just created replica use the following command and change head-osd to the one listed by `xtfs_repl -l`:



```
$> xtfs_repl -r head-OSD local-path-of-file
```

For further options see:

```
$> xtfs_repl -h
```

### 5.3.6 Automatic On-Close Replication

XtreemFS allows files to be automatically made read-only and replicated when being closed after having been written. On-close replication can e.g. be used to ensure that any file created on a volume is automatically replicated across multiple data centers. To configure the behavior of the on-close replication, `xtfs_repl` is used.

The number of replicas to be created when a file is closed can be set as follows:

```
$> xtfs_repl --ocr_factor_set 2
```

This will create two replicas of a file when the file is closed. Note that by setting the replication factor to 1 (default value), on-close replication will be switched off, which means that the file won't be replicated and will remain writable after having been closed. The current replication factor of a volume can be retrieved as follows:

```
$> xtfs_repl --ocr_factor_get
```

Moreover, it is possible to specify whether an automatically created replica will be filled immediately or on demand. By default, replicas will be filled on demand. This can be changed as follows:

```
$> xtfs_repl --ocr_full_set true
```

Thus, replicas will be filled in the background when the file is closed after having been written. To show whether replicas are automatically filled or not, execute the following command:

```
$> xtfs_repl --ocr_full_get
```



## Chapter 6

# Troubleshooting and Support

### 6.1 Logfiles

The logfiles for the XtreamFS services are located in `/var/log/xtreemfs`. The client generates no output, unless the `-f` and `-d INFO` options are specified.

### 6.2 Support

Please visit the [XtreamFS website at www.XtreamFS.org](http://www.XtreamFS.org) for links to the user mailing list and IRC channel.

### 6.3 Troubleshooting

**Problem:** The client hangs when opening/copying/creating a file but operations like `ls` or `mkdir` work.

**Solution:** This problem can occur when an OSD uses a UUID which resolves to an address that the client cannot (correctly) resolve. If you use e.g. `localhost:32640` as the UUID for the OSD, the client will try to contact the local machine instead of the machine on which the OSD runs. Check the status page of your Directory Service and check the UUID of the OSDs.

**Problem:** `xtfs_mount` does not print an error message but the volume is not mounted (i.e. not listed in the output of `mount`).

**Solution:** The client `xtfs_mount` automatically goes into background and does not print any error messages or warnings. Use the `-f` flag when mounting to prevent the client from going into background. All error messages will be printed to the console.



# Appendix A

## XtreemOS Integration

### XtreemFS Security Preparations

XtreemFS can be integrated in an existing XtreemOS VO security infrastructure. XtreemOS uses X.509 certificates to authenticate users in a Grid system, so the general setup is similar to a normal SSL-based configuration.

Thus, in an XtreemOS environment, certificates have to be created for the services as a first step. This is done by issuing a *Certificate Signing Request (CSR)* to the RCA server by means of the `create-server-csr` command. For further details, see the Section Using the RCA in the XtreemOS User Guide.

Signed certificates and keys generated by the RCA infrastructure are stored locally in PEM format. Since XtreemFS services are currently not capable of processing PEM certificates, keys and certificates have to be converted to PKCS12 and Java Keystore format, respectively.

Each XtreemFS service needs a certificate and a private key in order to be run. Once they have created and signed, the conversion has to take place. Assuming that certificate/private key pairs reside in the current working directory for the Directory Service, an MRC and an OSD (`ds.pem`, `ds.key`, `mrc.pem`, `mrc.key`, `osd.pem` and `osd.key`), the conversion can be initiated with the following commands:

```
$> openssl pkcs12 -export -in ds.pem -inkey ds.key \
    -out ds.p12 -name "DS"
$> openssl pkcs12 -export -in mrc.pem -inkey mrc.key \
    -out mrc.p12 -name "MRC"
$> openssl pkcs12 -export -in osd.pem -inkey osd.key \
    -out osd.p12 -name "OSD"
```

This will create three PKCS12 files (`ds.p12`, `mrc.p12` and `osd.p12`), each containing the private key and certificate for the respective service.

XtreemFS services need a *trust store* that contains all trusted Certification Authority certificates. Since all certificates created via the RCA have been signed by the XtreemOS CA, the XtreemOS CA certificate has to be included in the trust store. To create a new trust store containing the XtreemOS CA certificate, execute the following command:

```
$> keytool -import -alias xosrootca -keystore xosrootca.jks \
        -trustcacerts -file \
        /etc/xos/truststore/xtreemosrootcacert.pem
```

This will create a new Java Keystore `xosrootca.jks` with the XtreamOS CA certificate in the current working directory. The password chosen when asked will later have to be added as a property in the service configuration files.

Once all keys and certificates have been converted, the resulting files should be moved to `/etc/xos/xtreemfs/truststore/certs` as root:

```
# mv ds.p12 /etc/xos/xtreemfs/truststore/certs
# mv mrc.p12 /etc/xos/xtreemfs/truststore/certs
# mv osd.p12 /etc/xos/xtreemfs/truststore/certs
# mv xosrootca.jks /etc/xos/xtreemfs/truststore/certs
```

For setting up a *secured* XtreamFS infrastructure, each service provides the following properties:

```
# specify whether SSL is required
ssl.enabled = true

# server credentials for SSL handshakes
ssl.service_creds = /etc/xos/xtreemfs/truststore/certs/\
service.p12
ssl.service_creds.pw = xtreemfs
ssl.service_creds.container = pkcs12

# trusted certificates for SSL handshakes
ssl.trusted_certs = /etc/xos/xtreemfs/truststore/certs/\
xosrootca.jks
ssl.trusted_certs.pw = xtreemfs
ssl.trusted_certs.container = jks
```

`service.p12` refers to the converted file containing the credentials of the respective service. Make sure that all paths and passphrases (xtreemfs in this example) are correct.

## Appendix B

# Command Line Utilities

**xtfs\_cleanup** Deletes orphaned objects on an OSD or creates new metadata objects for orphaned files.

**xtfs\_lsvol** Lists the volumes on an MRC.

**xtfs\_mkvol** Creates a new volume on an MRC.

**xtfs\_mount** The XtreamFS client which mounts an XtreamFS volume locally on a machine.

**xtfs\_mrddbtool** Dumps an XML representation of the MRC database to a given directory in the MRC's local file system.

**xtfs\_repl** Displays and controls the replication of files and the replication policies.

**xtfs\_rmvol** Deletes a volume and all files on that volume from the MRC and the OSDs.

**xtfs\_stat** Displays XtreamFS specific file information such as the striping policy and the OSDs.

**xtfs\_sp** Displays and modifies the striping policy for a file, or the default striping policy for directories and volumes.

**xtfs\_scrub** Examines all files in a volume for incorrect file sizes and checksums. In case of incorrect file sizes, file sizes are corrected at the MRC.

**xtfs\_umount** Un-mounts a mounted XtreamFS volume.

# Index

- Access Policy, 6
  - Authorize All, 6
  - POSIX ACLs, 6
  - POSIX Permissions, 6
  - Volume ACLs, 6
- allow\_others option, 26
- allow\_root option, 26
- Architecture, 2
- Authentication, 3
- Authentication Provider, 11
  - NullAuthProvider, 11
  - SimpleX509AuthProvider, 11
  - XOAuthProvider, 12
- Authorization, 3
- Authorize All Access Policy, 6
- CA
  - Certificate Authority, 21
- Certificate, 3, 20
- Certificate Authority, 21
- Client, 3
- Create Volume, 23
- Credentials, 20
- Delete Volume, 24
- DIR, 2
- Directory Service, 2
- fileID, 30
- FUSE, 3
- init.d, 23
- Java KeyStore, 21
- JKS, 21
- Metadata, 2
- Metadata and Replica Catalog, 2
- Metadata Server, 2
- Mount, 26
- Mounting, 3
- MRC, 2
- NullAuthProvider, 11
- Object, 2
- Object Storage Device, 2
- Object-base File System, 2
- On-close Replication, 33
- OSD, 2
- OSD Selection Policy, 3
- PKCS#12, 20
- Policy
  - Access Policy, 6
  - OSD Selection Policy, 3
  - Striping Policy, 6
- POSIX ACLs Access Policy, 6
- POSIX Permissions Access Policy, 6
- RAIDo, 6
- Read-only Replication, 32
- Replication, 32, 33
  - on-close, 33
  - read-only, 32
- SimpleX509AuthProvider, 11
- SSL, 3
- Status Page, 23
- Storage Server, 2
- Stripe Size, 6
- Striping, 6
  - Stripe Size, 6
  - Striping Width, 6
- Striping Policy, 6
- Striping Width, 6
- Unmount, 26
- user\_allow\_other option, 26
- UUID, 10
- VFS, 3
- Volume, 2, 3
  - Create, 23
  - Delete, 24



- Mount, 26
- Un-mount, 26
- Volume ACLs Access Policy, 6
- Width, Striping Width, 6
- X.509, 3, 20
- XOSAuthProvider, 12
- xtfs\_mkvol, 23
- xtfs\_mount, 26
- xtfs\_repl, 32
- xtfs\_rmvol, 24
- xtfs\_sp, 30
- xtfs\_stat, 29
- xtfs\_umount, 26
- XtreemFS stat, 29
- XtreemFS striping policy tool, 30
- XtreemOS
  - Integration, 37
  - XtreemOS Certificates, 12