



DE LA RECHERCHE À L'INDUSTRIE

cea

www.cea.fr

Lecture on Parallel Filesystems

Lustre (4/5)

Jacques-Charles Lafoucriere

ENSIE| 2018

Unix Right



Unix User

Unix user identification

- User ID, associated to a User name
- Primary Group ID, associated to a Group name
- Additional group IDs, associated to Group names

```
$ id
uid=1007(lafoucrierejc) gid=1000(teachers) groups=1000(teachers),2000(guests),2001(hpcvm)
```



Unix Standard File access

3 components are associated to a file/directory

- User
 - Match a UID
 - File/Directory owner
- Group
 - Match a GID
 - File/Directory group
- Other
 - All others User = all users except owner

3 modes are associated to each component

- Read
- Write
- Execute

DE LA RECHERCHE À L'INDUSTRIE
cea

Usual Representation of Unix Rights

Diagram illustrating the usual representation of Unix rights (permissions) for three categories: User, Group, and Other. Each category has three permissions: Read (r), Write (w), and Execute (x). The diagram shows the bit set (solid line) and bit not set (dashed line) for each permission.

Bit Set

Bit Not Set

User Group Other

SFP 2018 Parallel Filesystems / Lustre | PAGE 5

DE LA RECHERCHE À L'INDUSTRIE
cea

Right Usage

File

- r access file content
- w modify file content
- x execute file
 - To run a shell script : need rx

Directory

- r list content
- w create, remove, rename an entry
- x cross the directory for searching
 - A file can be used even if the path cannot be listed

SFP 2018 Parallel Filesystems / Lustre | PAGE 6



Special Bits

Set UID

- On file: at execution change the effective user id by file owner id
- On directory: when an entry is created, use directory user in place of user id (not everywhere and dangerous)
- Display as s if x is on or S if x is off

Set GID

- On file: at execution change effective group id by file group id
- On directory: when an entry is created, use directory group in place of user group
- Display as s if x is on or S if x is off

Sticky bit

- On directory forbid remove of user not owned entries
- On file "stick" the execution code in swap device (no more used)

SFP 2018

Parallel Filesystems / Lustre | PAGE 7



How to Set/Change Unix Rights

chmod command line

Right rules

`[ugoa]*([-+=]([rwxXst]*|[ugo]))+([-+=])[0-7]+`

SFP 2018

Parallel Filesystems / Lustre | PAGE 8



X Chmod Flag Use

■ X: give x right if another x bit already exist

```
$ chmod 644 file ; ls -l file
-rw-r--r-- 1 LAFOUCRIEREJC None 0 May  6 15:13 file
$ chmod o+X file ; ls -l file
-rw-r--r-- 1 LAFOUCRIEREJC None 0 May  6 15:13 file
$ chmod u+x file ; ls -l file
-rwxr--r-- 1 LAFOUCRIEREJC None 0 May  6 15:13 file
$ chmod o+X file ; ls -l file
-rwxr--r-x 1 LAFOUCRIEREJC None 0 May  6 15:13 file
```

Access Control List



Context

Origin

- Standard Unix (POSIX) access rights are too restricted
- Need to be able to give/remove
 - Rights to a specific user/group
 - Independently of file/directory owner/group

Solution: ACL

- A set of data that informs an operating system about permissions or access rights that each user or group has to specific system objects such as directories or files
- A unique security attribute
- Defined as a POSIX standard

SFP 2018

Parallel Filesystems / Lustre | PAGE 11



ACL Definition

3 classes

- Owner
- Group
- Other

- Each ACL entry is
 - A class
 - A member of the class (no member == all members)
 - A set of rights rwx

SFP 2018

Parallel Filesystems / Lustre | PAGE 12



ACL Display

Example of ACL

```
# file: spoo
# owner: root
# group: root
user::rw-
user:bin:rw-
group::r--
mask::rw-
other::r--
```



ACL Display (Cont.)

A more complex of ACL

```
# file: .
# owner: root
# group: root
user::rwx
group::r-x
mask::rwx
other::r-x
default:user::rwx
default:group::rw-
default:group:bin:r--
default:mask::rw-
default:other::r-x
```



ACL Rules

Minimal ACL

- user::--- group::--- other::--- : match std Unix rights

Extended ACL

- Named entries
 - user:NAME:---
 - group:NAME:---
- Mask entry
 - mask:---
 - Mask named user and group entries

SFP 2018

Parallel Filesystems / Lustre | PAGE 15



ACL Rules (Cont.)

Default ACL

- default:user::---
- default:group::---
- default:user:NAME:---
- default:group:NAME:---
- default::other::---
- default::mask:---
- No role in access check
- Used to set ACL inherited at object creation
- Umask is ignored if an default ACL is present

SFP 2018

Parallel Filesystems / Lustre | PAGE 16



ACL Access Check Algorithm

ACL entries are looked at in the following order

- Owner
- Named users
- Owning or named groups
- Others

- Only a single entry determines access

SFP 2018

Parallel Filesystems / Lustre | PAGE 17



ACL Usage

How to know if an ACL exist?

```
$ ls -l file  
-rw-rw-r--+ 1 root root 0 May  6 16:45 file
```

How create/change/read ACL?

```
$ getfacl  
$ setfacl
```

SFP 2018

Parallel Filesystems / Lustre | PAGE 18



Read ACL

```
$ getfacl object
# file: object
# owner: root
# group: root
user::rw-
group::r--
other::r--
$ getfacl -c object
user::rw-
group::r--
other::r--
$
```



Modify ACL

With commands option

```
$ setfacl -m rule1,rule2,... object
$ setfacl -b object # remove all ACL entries
```

From an ACL file

```
$ getfacl -c object > template.acl # get a template
$ vi template.acl # edit ACL file
$ setfacl -M template.acl object
```



How to Setup ACL on Lustre

Activate ACL support on MDT

```
# mkfs -t lustre --mdt [...] --mountfsoption=user_xattr,acl /dev/MDT-Device
```

To check after MDT start

```
# lctl get_param -n mdc.test-MDT0000-mdc-*.connect_flags | grep acl  
acl
```

Extended Attributes



Extended Attribute

Definition

- name-value pairs associated with a file or directory on a filesystem

Applications

- POSIX ACLs
- Lustre metadata
 - Layout
 - FID
 - Parent FID
 - ...



Lustre Internal Extended Attribute

- trusted.som: size on mdt (not used)
- trusted.lov: file layout
- trusted.lma: obj FID (on MDT obj)
- trusted.lmv: directory layout
- trusted.dmv: defined to use xattr interface
- trusted.link: hard links vector (parent FID + others)
- trusted.fid: obj FID (on OST obj)
- trusted.version: defined to use xattr interface
- trusted.hsm: HSM EA
- trusted.lfsck_bitmap: lfsck namespace
- trusted.dummy: artefact



Root Squash

Objective

- Security feature
- Restrict super-user access rights from clients
 - Some client root is not root in Lustre namespace
- Useful if all clients does not share the same admin population

Implementation

- Root UID/GID is mapped to a defined UID/GID on the MGS

SFP 2018

Parallel Filesystems / Lustre | PAGE 27



Root Squash Configuration

Parameters

- `mdt.root_squash=UID:GID`
- `mdt.nosquash_nids="nids, nids, ..."`
 - Support use of jokers like `*@tcp`
 - List cleared with use of "clear" for the list

SFP 2018

Parallel Filesystems / Lustre | PAGE 28



Root Squash Configuration (Cont.)

Set root squash on MGS

```
# lctl conf_param test.mdt.root_squash=501:501
# lctl set_param mdt.test-MDT0000.root_squash=501:501
```

Test it

```
# touch file
# ls -l file
-rw-r--r-- 1 501 501 0 May  6 23:11 file
#
```

Thank you for your
attention



Lustre Glossary

| | | |
|------------------------------------------------------------------------------------------------------------------------------------------|--------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------------------------------------------------------|
| DNE - Distributed Namespace Environment - feature to aggregate multiple MDTs (possibly on many MDS's) into a single filesystem namespace | MDC - MetaData Client - client software layer that interfaces to the MDS | OSC - Object Storage Client - client software layer that interfaces to the OST |
| IDIF - OST object ID In FID - specific FID range reserved for compatibility with pre-DNE OST objects | MDD - Metadata Device Driver - MDS software layer that understands POSIX semantics for file access | OSD - Object Storage Device - server software layer that abstracts MDD and OFD access to underlying disk filesystems like Idiskfs and ZFS |
| IGIF - Inode and Generation In FID - specific FID range reserved for compatibility from Lustre 1.x MDT inode objects | MDS - MetaData Server - software service that manages access to filesystem namespace (inodes, paths, permission) requests from the client. | OSP - Object Storage Proxy - server software layer that interfaces from one MDS to the OSD on another MDS or another OSS |
| FID - File Identifier - unique 128-bit identifier for every object within a single filesystem. | MDT - MetaData Target - storage device that holds the filesystem metadata (attributes, inodes, directories, xattrs, etc) | OSS - Object Storage Server - software service that manages access to filesystem data (read, write, truncate, etc) |
| LMV - Logical Metadata Volume - client software layer that handles client (llite) access to multiple MDTs | MGS - Management Server - service that helps clients and servers with configuration | OST - Object Storage Target - storage device that holds the filesystem data (regular data files, not directories, xattrs, or other metadata) |
| LOD - Logical Object Device - MDS software layer that handles access to multiple MDTs and multiple OSTs | MGT - Management Target - storage device that holds the configuration logs | |
| LOV - Logical Object Volume - client software layer that handles client (llite) access to multiple OSTs | OFD - Object Filter Device - OSS software layer that handles file IO | |