

**Milestone Completion for
Implementation Milestone 1 of the
Distributed Namespace Project of
contract SFS-DEV-001.**

Revision History

Date	Revision	Author
Date (fixed)	Original	R. Henwood

Contents

Introduction.....	4
Subproject Description.....	4
Milestone Completion Criteria.....	4
Sanity.sh passed.....	4
mdsrate-create passed.....	4
Regression tests implemented and passed.....	5
Upgrade demonstrated.....	6
Downgrade demonstrated.....	8
Conclusion.....	8
Appendix A: sanity.sh screenshot.....	9
Appendix B: Downgrade demonstration.....	10

Introduction

The following milestone completion document applies to Subproject 2.1 - Remote Directories subproject within the OpenSFS Lustre Development contract SFS-DEV-001 signed 7/30/2011.

Subproject Description

Per the contract, DNE 1: Remote Directories is described as follows:

This subproject distributes the Lustre namespace over multiple metadata targets (MDTs) under administrative control using a Lustre-specific mkdir command. Whereas normal users are only able to create child directories and files on the same MDT as the parent directory, administrators can use this command to create a directory on a different MDT. The contents of any directory remain limited to a single MDT. Rename and hardlink operations between files and directories on different MDTs return EXDEV, forcing applications and utilities to treat them as if they are on different file systems. This limits the complexity of the implementation of this subproject while delivering capacity and performance scaling benefits for the entire namespace in aggregate.

Metadata update operations that span multiple MDTs are sequenced and synchronized to create and/or increment the link count on a MDT object before it is referenced by the remote directory entry and to update the remote directory entry before decrementing the link count and/or destroying the MDT object it referenced. Although this may result in an orphan MDT object under some failure conditions, it ensures that the Lustre namespace remains intact under any and all failure scenarios. All the other metadata operations avoid synchronous I/O and execute with full performance.

This project includes the implementation of OST FIDs (File Identifiers). These are required to overcome a limitation in the current 2.x Lustre protocol that would otherwise prevent a single file system from having more than 8 MDTs. Addressing this technical debt in the first subproject of DNE avoids protocol compatibility issues that would arise if this feature were implemented after Remote Directories were used in production.

Milestone Completion Criteria

Per the contract, three Implementation milestones have been defined by the Whamcloud. This document is concerned with completing the first Implementation milestone which is agreed as: "Demonstrate working DNE code. The sanity.sh and mdsrate-create tests will pass in a DNE environment. Suitable new regression tests for the remote directory functionality will be added and passed, including

functional Use Cases for upgrade and downgrade.” These requirements are enumerated as:

Demonstrate working DNE code

- sanity.sh passed.
- mdsrate-create passed.
- Regression tests implemented and passed.
- Upgrade demonstrated.
- Downgrade demonstrated.

These requirements are demonstrated below.

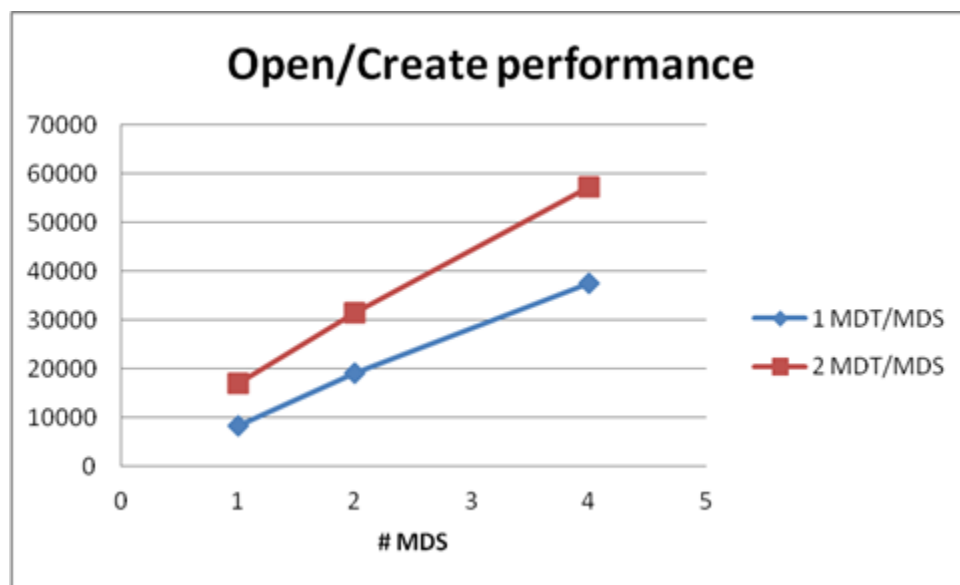
Sanity.sh passed

The results from the sanity.sh run are recorded in maloo at:

https://maloo.whamcloud.com/test_sets/4c945302-9bf0-11e1-8837-52540035b04c

A screenshot of the page is available as Appendix A: sanity.sh screenshot.

mdsrate-create passed



Data was collected from a Hyperion test run between April 12th and April 17th 2012. The test configuration included 100 clients, 4 threads on each client, each thread with an individual mount point. Each thread performs 10000 file open/create within unique directories. The units of the y-axis are completed file operations per second.

Regression tests implemented and passed

test_230 has been created to test DNE Remote Directories functionality:

```
test_230() {
    [ "$MDT_COUNT" -lt "2" ] && skip_env "skipping remote directory test" && return
    local MDTIDX=1

    mkdir -p $DIR/$tdir/test_230_local
    local mdt_idx=$(GETSTRIPE -M $DIR/$tdir/test_230_local)
    [ $mdt_idx -ne 0 ] &&
        error "create local directory on wrong MDT $mdt_idx"

    $LFS setdirstripe -i $MDTIDX $DIR/$tdir/test_230 ||
        error "create remote directory failed"
    local mdt_idx=$(GETSTRIPE -M $DIR/$tdir/test_230)
    [ $mdt_idx -ne $MDTIDX ] &&
        error "create remote directory on wrong MDT $mdt_idx"

    createmany -o $DIR/$tdir/test_230/t- 10 || error "create files on remote directory
failed"
    mdt_idx=$(GETSTRIPE -M $DIR/$tdir/test_230/t-0)
    [ $mdt_idx -ne $MDTIDX ] && error "create files on wrong MDT $mdt_idx"
    rm -r $DIR/$tdir || error "unlink remote directory failed"
}
run_test 230 "Create remote directory and files under the remote directory"
```

Upgrade demonstrated

Upgrade functionality is demonstrated as part of the test test_32c in conf-sanity:

https://maloo.whamcloud.com/test_sets/5cfc5278-9d2e-11e1-8587-52540035b04c

```
== conf-sanity test 32c: Upgrade with writeconf ===== 14:09:40 (1336932580)
Loading modules from /work/orion_release/orion/lustre-dev/lustre
../libcfs/libcfs/libcfs options: 'libcfs_panic_on_lbug=0'
debug=-1
subsystem_debug=all -lnet -lnd -pinger
../lnet/lnet/lnet options: 'networks=tcp accept=all'
gss/krb5 is not supported
/work/orion_release/orion/lustre-dev/lustre/Utils/tunefs.lustre
arch
commit
kernel
list
mdt
ost
shasums
Upgrading from disk2_2-ldiskfs.tar.bz2, created with:
  Commit: 2.2
  Kernel: 2.6.32-220.el6_lustre.g4554b65.x86_64
  Arch: x86_64
debug=-1
mount old MDT ....
mkfs new MDT....
mkfs.lustre: Warning: default mount option `errors=remount-ro' is missing
mount new MDT....
Mount client with 2 MDTs
Create the local directory and files on the old MDT
total: 10 creates in 0.04 seconds: 230.77 creates/second
Verify the MDT index of these files...Pass
Create the remote directory and files on new MDT
total: 10 creates in 0.05 seconds: 188.40 creates/second
Verify the MDT index of these files...Pass
Skip b1_8 images before we have 1.8 compatibility
Skip b1_8 images before we have 1.8 compatibility
Resetting fail_loc on all nodes...done.
```

Downgrade demonstrated

Downgrade demonstration is recorded on: <http://jira.whamcloud.com/browse/OSF-102>

A detailed record of the demonstration is recorded in Appendix B.

Conclusion

Implementation phase 1 has been completed according to the agreed criteria.

Appendix B: Downgrade demonstration

```
[root@testnode1 tests]# sh -vx upgrade_downgrade.sh
#Upgrade tests
#go to lustre 2.3 make old ldiskfs
cd /work/lustre-2.3/lustre/tests
+ cd /work/lustre-2.3/lustre/tests
testnode=${testnode:-`hostname`}
hostname
++ hostname
+ testnode=testnode1
LOAD=y sh llmount.sh
+ LOAD=y
+ sh llmount.sh
Loading modules from /work/lustre-2.3/lustre/tests/..
debug=vfstrace rpctrace dlctrace neterror ha config ioctl super
subsystem debug=all -lnet -lnd -pinger
../lnet/lnet/lnet options: 'networks=tcp accept=all'
gss/krb5 is not supported
quota/lquota options: 'hash_lqs_cur_bits=3'
set -vx
+ set -vx
../utils/mkfs.lustre --reformat --mgs --mdt --device-size=1048576 /tmp/lustre-mdt-2.3
+ ../utils/mkfs.lustre --reformat --mgs --mdt --device-size=1048576 /tmp/lustre-mdt-2.3
Permanent disk data:
Target: lustre-MDTffff
Index: unassigned
Lustre FS: lustre
Mount type: ldiskfs
Flags: 0x75
(MDT MGS needs_index first_time update )
Persistent mount opts: user_xattr,errors=remount-ro
Parameters:
formatting backing filesystem ldiskfs on /dev/loop0
target name lustre-MDTffff
4k blocks 262144
options -I 512 -i 2048 -q -O dirdata,uninit_bg,dir_nlink,huge_file,flex_bg -E
lazy_journal_init -F
mkfs_cmd = mke2fs -j -b 4096 -L lustre-MDTffff -I 512 -i 2048 -q -O
dirdata,uninit_bg,dir_nlink,huge_file,flex_bg -E lazy_journal_init -F /dev/loop0 262144
Writing CONFIGS/mountdata
../utils/mkfs.lustre --reformat --mgsnode=$testnode --ost --device-size=1048576
/tmp/lustre-ost-2.3
+ ../utils/mkfs.lustre --reformat --mgsnode=testnode1 --ost --device-size=1048576
/tmp/lustre-ost-2.3
Permanent disk data:
Target: lustre-OSTffff
Index: unassigned
Lustre FS: lustre
Mount type: ldiskfs
Flags: 0x72
(OST needs_index first_time update )
Persistent mount opts: errors=remount-ro,extents,mballoc
Parameters: mgsnode=192.168.122.162@tcp
formatting backing filesystem ldiskfs on /dev/loop0
target name lustre-OSTffff
4k blocks 262144
options -I 256 -q -O extents,uninit_bg,dir_nlink,huge_file,flex_bg -G 256 -E
resize=4290772992,lazy_journal_init -F
mkfs_cmd = mke2fs -j -b 4096 -L lustre-OSTffff -I 256 -q -O
extents,uninit_bg,dir_nlink,huge_file,flex_bg -G 256 -E
resize=4290772992,lazy_journal_init -F /dev/loop0 262144
Writing CONFIGS/mountdata
mount -t lustre -o loop,user_xattr,acl /tmp/lustre-mdt-2.3 /mnt/mds
+ mount -t lustre -o loop,user_xattr,acl /tmp/lustre-mdt-2.3 /mnt/mds
mount -t lustre -o loop /tmp/lustre-ost-2.3 /mnt/ost1
+ mount -t lustre -o loop /tmp/lustre-ost-2.3 /mnt/ost1
```



```

mount -t lustre $testnode:/lustre /mnt/lustre
+ mount -t lustre testnode1:/lustre /mnt/lustre
cp /etc/fstab /mnt/lustre
+ cp /etc/fstab /mnt/lustre
cp /etc/hosts /mnt/lustre
+ cp /etc/hosts /mnt/lustre
umount /mnt/lustre
+ umount /mnt/lustre
umount /mnt/ost1
+ umount /mnt/ost1
umount /mnt/mds
+ umount /mnt/mds
losetup -d /dev/loop0
+ losetup -d /dev/loop0
losetup -d /dev/loop1
+ losetup -d /dev/loop1
losetup -d /dev/loop2
+ losetup -d /dev/loop2
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop3
+ losetup -d /dev/loop3
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop4
+ losetup -d /dev/loop4
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop5
+ losetup -d /dev/loop5
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop6
+ losetup -d /dev/loop6
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop7
+ losetup -d /dev/loop7
ioctl: LOOP_CLR_FD: No such device or address
LOAD=y sh llmountcleanup.sh
+ LOAD=y
+ sh llmountcleanup.sh
Stopping clients: testnode1 /mnt/lustre (opts:-f)
Stopping clients: testnode1 /mnt/lustre2 (opts:-f)
osd_ldiskfs 296768 0
fsfilt_ldiskfs 119600 0
mdd 426496 3 osd_ldiskfs,cmm,mdt
ldiskfs 354264 2 osd_ldiskfs,fsfilt_ldiskfs
jbd2 101384 3 osd_ldiskfs,fsfilt_ldiskfs,ldiskfs
crc16 35328 1 ldiskfs
obdclass 1109104 29
llite_lloop,lustre,obdfilter,ost,osd_ldiskfs,cmm,fsfilt_ldiskfs,mdt,mdd,mds,mgs,mgc,lov,
osc,mdc,lmv,fid,fld,lquota,ptlrpc
lvfs 72256 22
llite_lloop,lustre,obdfilter,ost,osd_ldiskfs,cmm,fsfilt_ldiskfs,mdt,mdd,mds,mgs,mgc,lov,
osc,mdc,lmv,fid,fld,lquota,ptlrpc,obdclass
libcfs 344320 24
llite_lloop,lustre,obdfilter,ost,osd_ldiskfs,cmm,fsfilt_ldiskfs,mdt,mdd,mds,mgs,mgc,lov,
osc,mdc,lmv,fid,fld,lquota,ptlrpc,obdclass,lvfs,ksocklnd,lnet
exportfs 39296 2 fsfilt_ldiskfs,nfsd
modules unloaded.
echo "go to DNE branch do upgrade"
+ echo 'go to DNE branch do upgrade'
go to DNE branch do upgrade
cd /work/lustre-dne/lustre/tests
+ cd /work/lustre-dne/lustre/tests
LOAD=y sh llmount.sh
+ LOAD=y
+ sh llmount.sh
Loading modules from /work/lustre-dne/lustre/tests/..
../libcfs/libcfs/libcfs options: 'libcfs_panic_on_lbug=0'
debug=vfstrace rpctrace dlmtrace neterror ha config ioctl super
subsystem_debug=all -lnet -lnd -pinger
../lnet/lnet/lnet options: 'networks=tcp accept=all'
gss/krb5 is not supported
mount -t lustre -o loop,user_xattr,acl,abort_recov,write_conf /tmp/lustre-mdt-2.3
/mnt/mds1

```

```

+ mount -t lustre -o loop,user_xattr,acl,abort_recov,write_conf /tmp/lustre-mdt-2.3
/mnt/mds1
../utils/mkfs.lustre --reformat --mgsnode=$testnode --mdt --device-size=1048576 --index
1 /tmp/lustre-mdt-new
+ ../utils/mkfs.lustre --reformat --mgsnode=testnode1 --mdt --device-size=1048576
--index 1 /tmp/lustre-mdt-new
Permanent disk data:
Target: lustre:MDT0001
Index: 1
Lustre FS: lustre
Mount type: ldiskfs
Flags: 0x61
(MDT first_time update )
Persistent mount opts: user_xattr,errors=remount-ro
Parameters: mgsnode=192.168.122.162@tcp
formatting backing filesystem ldiskfs on /dev/loop1
target name lustre:MDT0001
4k blocks 262144
options -I 512 -i 2048 -q -O dirdata,uninit_bg,dir_nlink,huge_file,flex_bg -E
lazy_journal_init -F
mkfs_cmd = mke2fs -j -b 4096 -L lustre:MDT0001 -I 512 -i 2048 -q -O
dirdata,uninit_bg,dir_nlink,huge_file,flex_bg -E lazy_journal_init -F /dev/loop1 262144
Writing CONFIGS/mountdata
mount -t lustre -o loop,user_xattr,acl,abort_recov,write_conf /tmp/lustre-mdt-new
/mnt/mds2
+ mount -t lustre -o loop,user_xattr,acl,abort_recov,write_conf /tmp/lustre-mdt-new
/mnt/mds2
mount -t lustre -o loop,abort_recov,write_conf /tmp/lustre-ost-2.3 /mnt/ost1
+ mount -t lustre -o loop,abort_recov,write_conf /tmp/lustre-ost-2.3 /mnt/ost1
mount -t lustre $testnode:/lustre /mnt/lustre
+ mount -t lustre testnode1:/lustre /mnt/lustre
diff /mnt/lustre/fstab /etc/fstab || { echo "the file is diff1" && exit 1; }
+ diff /mnt/lustre/fstab /etc/fstab
diff /mnt/lustre/hosts /etc/hosts || { echo "the file is diff1" && exit 1; }
+ diff /mnt/lustre/hosts /etc/hosts
../utils/lfs setdirstripe -i 1 /mnt/lustre/test_mdt1 || { echo "create remote directory
failed" && exit 1; }
+ ../utils/lfs setdirstripe -i 1 /mnt/lustre/test_mdt1
mdt_idx=$(../utils/lfs getstripe -M /mnt/lustre/test_mdt1)
../utils/lfs getstripe -M /mnt/lustre/test_mdt1
++ ../utils/lfs getstripe -M /mnt/lustre/test_mdt1
+ mdt_idx=1
[ $mdt_idx -ne 1 ] && { echo "create remote directory on wrong MDT" && exit 1; }
+ '[' 1 -ne 1 ']'

mkdir /mnt/lustre/test_mdt1/dir
+ mkdir /mnt/lustre/test_mdt1/dir
mdt_idx=$(../utils/lfs getstripe -M /mnt/lustre/test_mdt1/dir)
../utils/lfs getstripe -M /mnt/lustre/test_mdt1/dir
++ ../utils/lfs getstripe -M /mnt/lustre/test_mdt1/dir
+ mdt_idx=1
[ $mdt_idx -ne 1 ] && { echo "create remote directory on wrong MDT" && exit 1; }
+ '[' 1 -ne 1 ']'

cp /mnt/lustre/hosts /mnt/lustre/test_mdt1/dir/hosts
+ cp /mnt/lustre/hosts /mnt/lustre/test_mdt1/dir/hosts
cp /mnt/lustre/fstab /mnt/lustre/test_mdt1/dir/fstab
+ cp /mnt/lustre/fstab /mnt/lustre/test_mdt1/dir/fstab
echo "downgrade DNE to single MDT"
+ echo 'downgrade DNE to single MDT'
downgrade DNE to single MDT
mkdir /mnt/lustre/test_mdt1_backup
+ mkdir /mnt/lustre/test_mdt1_backup
cp -R /mnt/lustre/test_mdt1/ /mnt/lustre/test_mdt1_backup/
+ cp -R /mnt/lustre/test_mdt1/ /mnt/lustre/test_mdt1_backup/
../utils/lctl dk > /tmp/debug.out
+ ../utils/lctl dk
umount /mnt/lustre/
+ umount /mnt/lustre/
umount /mnt/mds2
+ umount /mnt/mds2
umount /mnt/mds1
+ umount /mnt/mds1

```

```
umount /mnt/ost1
+ umount /mnt/ost1
LOAD=y sh llmountcleanup.sh
+ LOAD=y
+ sh llmountcleanup.sh
Stopping clients: testnode1 /mnt/lustre (opts:-f)
Stopping clients: testnode1 /mnt/lustre2 (opts:-f)
osd_ldiskfs 343056 0
fsfilt_ldiskfs 43776 0
ldiskfs 354392 2 osd_ldiskfs,fsfilt_ldiskfs
mdd 338320 3 osd_ldiskfs,cmm,mdt
fld 113776 7 osd_ldiskfs,lod,obdfilter,cmm,mdt,lmv,fid
obdclass 1221936 28
llite_loop,lustre,osd_ldiskfs,osp,lod,obdfilter,ost,cmm,mdt,mdd,mgs,mgc,lov,osc,mdc,lmv,
,fid,fld,ptlrpc
lvfs 59024 22
llite_loop,lustre,osd_ldiskfs,fsfilt_ldiskfs,osp,lod,obdfilter,ost,cmm,mdt,mdd,mgs,mgc,
lov,osc,mdc,lmv,fid,fld,ptlrpc,obdclass
libcfs 344192 24
llite_loop,lustre,osd_ldiskfs,fsfilt_ldiskfs,osp,lod,obdfilter,ost,cmm,mdt,mdd,mgs,mgc,
lov,osc,mdc,lmv,fid,fld,ptlrpc,obdclass,lvfs,ksocklnd,lnet
jbd2 101384 2 osd_ldiskfs,ldiskfs
crcl6 35328 1 ldiskfs
modules unloaded.
losetup -d /dev/loop0
+ losetup -d /dev/loop0
losetup -d /dev/loop1
+ losetup -d /dev/loop1
losetup -d /dev/loop2
+ losetup -d /dev/loop2
losetup -d /dev/loop3
+ losetup -d /dev/loop3
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop4
+ losetup -d /dev/loop4
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop5
+ losetup -d /dev/loop5
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop6
+ losetup -d /dev/loop6
ioctl: LOOP_CLR_FD: No such device or address
losetup -d /dev/loop7
+ losetup -d /dev/loop7
ioctl: LOOP_CLR_FD: No such device or address
cd /work/lustre-2.3/lustre/tests
+ cd /work/lustre-2.3/lustre/tests
LOAD=y sh llmount.sh
+ LOAD=y
+ sh llmount.sh
Loading modules from /work/lustre-2.3/lustre/tests/..
debug=vfstrace rpctrace dlmtrace neterror ha config ioctl super
subsystem_debug=all -lnet -lnd -pinger
../lnet/lnet/lnet options: 'networks=tcp accept=all'
gss/krb5 is not supported
quota/lquota options: 'hash_lqs_cur_bits=3'
/work/lustre-2.3/lustre/utlils/tunefs.lustre --writeconf /tmp/lustre-mdt-2.3
+ /work/lustre-2.3/lustre/utlils/tunefs.lustre --writeconf /tmp/lustre-mdt-2.3
checking for existing Lustre data: found CONFIGS/mountdata
Reading CONFIGS/mountdata
Read previous values:
Target: lustre-MDT0000
Index: 0
Lustre FS: lustre
Mount type: ldiskfs
Flags: 0x5
(MDT MGS )
Persistent mount opts: user_xattr,errors=remount-ro
Parameters:
Permanent disk data:
Target: lustre-MDT0000
Index: 0
```

```
Lustre FS: lustre
Mount type: ldiskfs
Flags: 0x105
(MDT MGS writeconf )
Persistent mount opts: user_xattr,errors=remount-ro
Parameters:
Writing CONFIGS/mountdata
/work/lustre-2.3/lustre/utils/tunefs.lustre --writeconf /tmp/lustre-ost-2.3
+ /work/lustre-2.3/lustre/utils/tunefs.lustre --writeconf /tmp/lustre-ost-2.3
checking for existing Lustre data: found CONFIGS/mountdata
Reading CONFIGS/mountdata
Read previous values:
Target: lustre-OST0000
Index: 0
Lustre FS: lustre
Mount type: ldiskfs
Flags: 0x2
(OST )
Persistent mount opts: errors=remount-ro,extents,malloc
Parameters: mgsnode=192.168.122.162@tcp
Permanent disk data:
Target: lustre-OST0000
Index: 0
Lustre FS: lustre
Mount type: ldiskfs
Flags: 0x102
(OST writeconf )
Persistent mount opts: errors=remount-ro,extents,malloc
Parameters: mgsnode=192.168.122.162@tcp
Writing CONFIGS/mountdata
mount -t lustre -o loop,user_xattr,acl,abort_recov /tmp/lustre-mdt-2.3 /mnt/mds
+ mount -t lustre -o loop,user_xattr,acl,abort_recov /tmp/lustre-mdt-2.3 /mnt/mds
mount -t lustre -o loop,abort_recov /tmp/lustre-ost-2.3 /mnt/ost1
+ mount -t lustre -o loop,abort_recov /tmp/lustre-ost-2.3 /mnt/ost1
mount -t lustre $testnode:/lustre /mnt/lustre
+ mount -t lustre testnode1:/lustre /mnt/lustre
diff /mnt/lustre/fstab /etc/fstab || { echo "the file is different" && exit 1; }
+ diff /mnt/lustre/fstab /etc/fstab
diff /mnt/lustre/hosts /etc/hosts || { echo "the file is different" && exit 1; }
+ diff /mnt/lustre/hosts /etc/hosts
diff /mnt/lustre/test_mdt1_backup/test_mdt1/dir/fstab /etc/fstab || { echo "the file is
different" && exit 1; }
+ diff /mnt/lustre/test_mdt1_backup/test_mdt1/dir/fstab /etc/fstab
diff /mnt/lustre/test_mdt1_backup/test_mdt1/dir/hosts /etc/hosts || { echo "the file is
different" && exit 1; }
+ diff /mnt/lustre/test_mdt1_backup/test_mdt1/dir/hosts /etc/hosts
umount /mnt/lustre
+ umount /mnt/lustre
umount /mnt/mds
+ umount /mnt/mds
umount /mnt/ost1
+ umount /mnt/ost1
```