

# Final Report for the MDT-OST Consistency Subproject of the Lustre File System Checker of the SFS-DEV-001 Contract

---

## Revision History

<b>Date</b>	<b>Revision</b>	<b>Author</b>
07/31/14	Original	R. Henwood
08/07/14	Improvements from OpenSFS team	R. Henwood

# Contents

- Executive Summary.....3
- Statement of Work.....3
- Summary of Solution Architecture.....4
- Acceptance Criteria.....5
- Summary of High Level Design.....6
- Summary of Implementation.....7
- Summary of Demonstration.....8
- Documentation.....8

## Executive Summary

This document finalizes the activities undertaken during the Lustre\* File System Checker, Sub Project 3.2: MDT-OST Consistency project within the OpenSFS Lustre Development contract SFS-DEV-001 signed July 30<sup>th</sup> 2011.

Notable highlights of this project include:

- Demonstrated scalable consistency checking over multiple MDTs tested on the OpenSFS Cluster.
- MDT-OST functionality was implemented with 12886 new lines of code and was completed and landed on Lustre Master for inclusion into Lustre release 2.6 on July 31<sup>st</sup> 2014.
- All assets generated for OpenSFS during the project are attached to the OpenSFS wiki: [http://wiki.opensfs.org/Contract\\_SFS-DEV-001](http://wiki.opensfs.org/Contract_SFS-DEV-001)

## Statement of Work

MDT-OST consistency implements functionality for distributed verification and repair of the MDT-inode-to-OST-object layout. The LFSC 2 functionality works in a distributed fashion with the MDT iterating over the inodes to check:

- the file layout (LOV EA) and verify the ostid therein
- the OST objects referenced by the file layout exist
- each OST object has a back reference to the correct MDT inode

Incorrect or missing back pointers on the OST objects are corrected, and missing objects recreated when detected.

The UID and GID of OST objects are also verified to match that of the MDT inode to ensure correct quota allocation. After the MDT iteration is complete, any unreferenced OST objects are linked into the Lustre .lustre/lost+found/ directory.

Subproject 3.1, 3.1.5 and this sub-project (3.2) altogether constitute a complete, scalable replacement of the existing e2fsprogs-based lfsck utility. This allows the distributed checking and repair of Lustre inter-server state for while the file system is

\*Other names and brands maybe the property of others.

on-line. LFSCCK supports multiple MDTs present on DNE file systems and will run simultaneously on multiple MDTs. MDT to MDT consistency is not supported in this phase of the project.

The complete scope statement was agreed on 2013-05-01 and is available at:

[http://wiki.opensfs.org/images/d/d3/LFSCCK\\_MDT-OSTConsistency\\_ScopeStatement.pdf](http://wiki.opensfs.org/images/d/d3/LFSCCK_MDT-OSTConsistency_ScopeStatement.pdf)

## Summary of Solution Architecture

Normal files on a Lustre\* file system (i.e. non-directory) are composed of one MDT object (the parent) and zero or more OST object(s) (or children). The parent resides on the MDT, and records the file layout information in the Logical Object Volume Extended Attribute (LOV EA) for the children belonging to the file. With the file layout information, a client can locate the specified OST object. To ensure data integrity, each child object on the related OST also records its parent identifier information to indicate to which file the OST object belongs. Under normal operation, the file layout information stored on the parent will be consistent with the parent identifier information stored in its children. With a production system, however, a error/failure condition may arise that can cause the parent-child pointers to become inconsistent. The inconsistency includes the following cases:

1. **Dangling reference.** *MDT-object1* claims *OST-object1* is its child, but on the OST, *OST-object1* does not exist, or is not initialized (and does not recognize *MDT-object1* as its parent).
2. **Unreferenced OST object.** *OST-object1* claims *MDT-object1* is its parent, but on the MDT, *MDT-object1* does not exist, or is not initialized (and does not recognize *OST-object1* as its child).
3. **Mismatched reference.** *MDT-object1* claims *OST-object1* is its child, but *OST-object1* claims that its parent is *MDT-object2*. On the MDT, *MDT-object2* doesn't exist, or it doesn't recognize *OST-object1* as its child.
4. **Multiple references.** *MDT-object1* claims *OST-object1* is its child, but *OST-object1* claims that its parent is *MDT-object2* rather than *MDT-object1*. *MDT-object2* recognizes *OST-object1* as its child.

A further type of inconsistency between an MDT and an OST is concerned with quota. Both MDT objects and OST objects have ownership information (UID and GID) to indicate which user the file/object belongs to. If the owner information for the MDT

object and OST object(s) belonging to the same file are inconsistent, then quota will be inaccurate.

To fix the MDT-OST inconsistencies identified above, a scan of the whole system is needed including both the MDT and OSTs is performed. The existing object-table based iteration, implemented in LFSCK Phase I, scans the whole system. In addition to MDT-OST inconsistencies, all previous LFSCK functionality from phases 1 and 1.5 will execute and correct the file system.

## Acceptance Criteria

The criteria that must be met for a successfully acceptance are:

1. The administrator can start and stop MDT-OST consistency check/repair through userspace commands.
2. The administrator can monitor MDT-OST consistency check/repair.
3. The administrator can resume MDT-OST consistency check/repair from the latest checkpoint.
4. The administrator can control the rate of scanning for MDT-OST consistency check/repair.
5. LFSCK will repair file for which the parent has a dangling reference.
6. LFSCK will repair an unreferenced OST object.
7. LFSCK will repair an unmatched referenced MDT object and OST object pair.
8. LFSCK will repair a repeat referenced OST object.
9. LFSCK will repair inconsistent file owner information.
10. The administrator can upgrade an OST from Lustre 1.8.
11. The Lustre system is available during the LFSCK for MDT-OST consistency check/repair.

The complete Solution Architecture including Acceptance Criteria was agreed on 2013-06-05 and is available at:

[http://wiki.opensfs.org/images/e/ea/LFSCK\\_MDT-OSTConsistency\\_SolutionArchitecture.pdf](http://wiki.opensfs.org/images/e/ea/LFSCK_MDT-OSTConsistency_SolutionArchitecture.pdf)

## Summary of High Level Design

MDT-OST consistency scanning (along with the other phases of LFSCK) is a highly complex implementation task. The complete High Level Design document weighs in at over 9000 words. That document describes in detail the following considerations:

- Identifying inconsistencies efficiently in a distributed environment.
- Tracing LFSCK during operation.
- Controlling and monitoring LFSCK from user space.
- Design of the independent LFSCK check/repair engines.
- Strategy to repair different inconsistencies in a distributed environment.
- Notification of events to the changelog.
- Recovery in the event of failure of any part of the system.
- Changes needed to the wire protocol to support MDT-OST consistency check/repair.
- Changes to the OSD API.
- Control and resolution of race conditions with LFSCK and other operations.
- On-disk layout changes.
- Interoperability and compatibility with other features.

The complete High Level Design was agreed on 2013-07-23 and is available at:

[http://wiki.opensfs.org/images/b/b9/LFSCK\\_MDT-OSTConsistency\\_HighLevelDesign.pdf](http://wiki.opensfs.org/images/b/b9/LFSCK_MDT-OSTConsistency_HighLevelDesign.pdf)

## Summary of Implementation

LFSCCK 2: MDT-OST consistency is implemented in the following patches:

ID	Description
<a href="#">591f15e</a>	<a href="#">LU-4106 scrub: Trigger OI scrub properly</a>
<a href="#">d628a95</a>	<a href="#">LU-3951 fsck: OST object inconsistency self detect/repair</a>
<a href="#">b3e6eda</a>	<a href="#">LU-3951 fsck: LWP connection from OST-x to MDT-y</a>
<a href="#">bf15de9</a>	<a href="#">LU-3950 fsck: control LFSCCK on all devices via single command</a>
<a href="#">80054e6</a>	<a href="#">LU-3594 fsck: repair inconsistent OST object owner</a>
<a href="#">cd21da7</a>	<a href="#">LU-3593 fsck: repair inconsistent layout EA</a>
<a href="#">aa87bcd</a>	<a href="#">LU-3592 fsck: repair multiple referenced OST object</a>
<a href="#">74c1059</a>	<a href="#">LU-3591 fsck: repair unmatched MDT-OST objects pairs</a>
<a href="#">aa1e2e7</a>	<a href="#">LU-3590 fsck: repair MDT object with dangling reference</a>
<a href="#">f792da4</a>	<a href="#">LU-3569 ofd: packing ost_idx in all IDIF</a>
<a href="#">2a271b4</a>	<a href="#">LU-3336 fsck: create new MDT object or exchange OST objects</a>
<a href="#">ac78f55</a>	<a href="#">LU-3336 fsck: recreate the lost MDT object</a>
<a href="#">416839c</a>	<a href="#">LU-3336 fsck: namespace visible lost+found directory</a>
<a href="#">c61c112</a>	<a href="#">LU-3336 fsck: regenerate lost layout EA</a>
<a href="#">56fdc3c</a>	<a href="#">LU-3336 fsck: orphan OST objects iteration (2)</a>
<a href="#">5b4f73b</a>	<a href="#">LU-3336 fsck: orphan OST objects iteration (1)</a>
<a href="#">4425fa1</a>	<a href="#">LU-3336 fsck: use rbtree to record OST object accessing</a>
<a href="#">3e09d63</a>	<a href="#">LU-3335 osd: use local transaction directly inside OSD</a>
<a href="#">cc581f3</a>	<a href="#">LU-1267 fsck: enhance API for MDT-OST consistency</a>
<a href="#">7a998f3</a>	<a href="#">LU-1267 fsck: enhance RPCs (2) for MDT-OST consistency</a>
<a href="#">a956e98</a>	<a href="#">LU-1267 fsck: enhance RPCs (1) for MDT-OST consistency</a>
<a href="#">edb1bd9</a>	<a href="#">LU-1267 fsck: framework (3) for MDT-OST consistency</a>
<a href="#">f946d82</a>	<a href="#">LU-1267 fsck: framework (2) for MDT-OST consistency</a>
<a href="#">141a375</a>	<a href="#">LU-1267 fsck: rebuild LAST_ID</a>
<a href="#">ff36f471</a>	<a href="#">LU-1267 fsck: framework (1) for MDT-OST consistency</a>

The Implementation milestone was agreed on 2013-09-27 and is available at:

[http://wiki.opensfs.org/images/e/ee/LFSCCK\\_MDT-OSTConsistency\\_Implementation.pdf](http://wiki.opensfs.org/images/e/ee/LFSCCK_MDT-OSTConsistency_Implementation.pdf)

## Summary of Demonstration

LFSCCK 2: Successfully completed the following demonstration criteria:

- Correctness: coded into sanity-lfscck.sh and sanity-scrub.sh
- Standard review tests: sanity, sanityn, replay-single, conf-sanity, recovery-small, replay-ost-single, insanity, sanity-quota, sanity-sec, lustre-rsync-test, lnet-selftest, and mmp
- Measure performance of LFSCCK 2 against a single MDT device without inconsistencies
- Measure performance of LFSCCK 2 against a single MDT device with inconsistencies
- Impact of LFSCCK 2 on small file create performance on a single MDT with inconsistencies
- Measure performance of LFSCCK 2 against multiple MDT devices during inconsistency repair.

Complete results were agreed on 2013-11-13 and are available at:

[http://wiki.opensfs.org/images/0/08/LFSCCK\\_MDT-OSTConsistency\\_Demonstration.pdf](http://wiki.opensfs.org/images/0/08/LFSCCK_MDT-OSTConsistency_Demonstration.pdf)

## Documentation

Documentation was completed as part of issue <http://review.whamcloud.com/9068> and landed on change <https://jira.hpdd.intel.com/browse/LUDOC-155>.

In addition, a design document for LFSCCK 2 has been included in the Lustre source code tree in the 'Documentation' folder.