



Whamcloud

Lustre 2.16 and Beyond

Andreas Dilger

Lustre Principal Architect

Planned Feature Release Highlights

▶ **2.16** opening to land new feature patches

- LNet IPv6 addressing – allow 160-bit NIDs, more flexible server configuration (SuSE)
- Optimized Directory Traversal (WBC1) – cross-directory stathread (WC)

▶ **2.17** has several major features already lined up

- Client-side data compression – use client CPU to reduce network and storage usage (WC)
- Metadata Writeback Cache (WBC2) – low latency file operations in client RAM (WC)
- File Level Redundancy - Erasure Coding (EC) – efficiently store file redundancy

▶ **2.18** feature proposals in early discussion stages

- Lustre Metadata Redundancy (LMR1) – MDT0000 service redundancy

LNet Improvements

(2.15/2.16)



▶ Multiple TCP sockets for 100GigE+ performance ([LU-12815](#), WC)

- Add `conns_per_peer=N` for `sock1nd` (4.1GB/s->**9.5GB/s** on 100GbE)
- Auto-configure based on interface speed (e.g. 10Gbps=>2, 100Gbps=>4, ...)

▶ LNet Network Selection Policy (UDSP) ([LU-9121](#), WC)

- Allow policies for local/remote interface prioritization by NID
 - e.g. primary IB with TCP backup, select "best" router NID for client/server

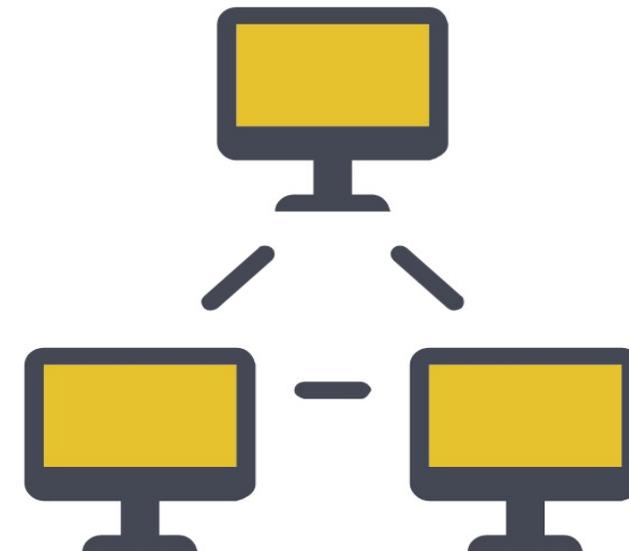
2.15

2.16 ▶ IPv6 NID support ([LU-10391](#), SuSE)

- Variable-sized NIDs (8-bit type, 8-bit size, 16-bit network, 128-bit+ address)
- Interoperable with existing current LNDs whenever possible

▶ Simplified/dynamic server node addressing ([LU-14668](#), WC)

- Detect added/changed server interfaces automatically ([LU-10360](#))
- Reduce (and eventually eliminate) static NIDs in Lustre config logs



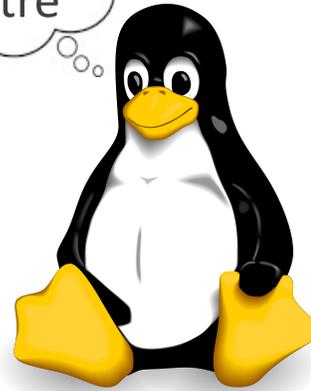
Client Improvements

▶ **GPU Direct RDMA** - directly into GPU, bypass CPU ([LU-14798](#), WC, NVIDIA, HPE)

- A100 2x200Gb IB **36GB/s** write, **39GB/s** read, **174GB/s** with 8x200Gb IB

▶ **Parallel large DIO** optimization ([LU-13798](#), [LU-13799](#), HPE, WC)

- Improve single-thread `read()/write()` (1.5GB/s->**15.8GB/s!**)
- Particular benefits for AIO/DIO and `io_uring` in client kernels 5.1 and later



2.15 ▶ Improved "`lfs find -printf`" option for scanning files ([LU-10378](#), ORNL)

2.16 ▶ `o2ib1nd` cleanups for in-kernel OFED ([LU-8874](#), ORNL)

▶ Buffered/DIO/mmap performance/efficiency improvements ([LU-13805](#), WC)

▶ Ongoing code style/structure cleanup for upstream submission (ORNL)

▶ Ongoing updates for newer kernels (ORNL, SuSE)

Backend OSD Improvements

- ▶ Parallel e2fsck for pass2/3 (directory entries, name linkage) ([LU-14679](#), WC)
 - Now slowest part of e2fsck (was 7% of total time, now **70%** after pass1/pass5 speedups)
- ▶ ZFS 2.1 dRAID VDEVs - declustered parity and hot space (LLNL, HPE, Intel)
- ▶ `fa1locate()` and `FALLOCATE_FL_PUNCH_HOLE` for ZFS ([LU-14157](#), AEON)
- ▶ Improved `ldiskfs` `mballoc` efficiency for large/full filesystems ([LU-14438](#), Google, WC)
 - $O(1)$ lookup of power-of-two free space, $O(\log N)$ lookup of other sizes
- ▶ Improved `ldiskfs` "`-o discard`" efficiency ([LU-14712](#), Kuaishou, WC)
 - Allow real-time TRIM of flash storage to maintain peak performance
- ▶ OST object directory scalability for large OSTs ([LU-11912](#), WC)
 - Large OSTs (500-1000TB) have billions of objects, only 32 dirs per MDT!
 - Wider dir fanout not better, object create/remove access all dirs randomly
 - New OST FID Sequences more often (e.g. 32M vs. 4B objs), retire old SEQ
 - Groups objects by age to limit directory size and improve efficiency

Batched Cross-Directory Statahead

(WC 2.16)



▶ **Batched RPCs** for multi-update operations ([LU-13045](#))

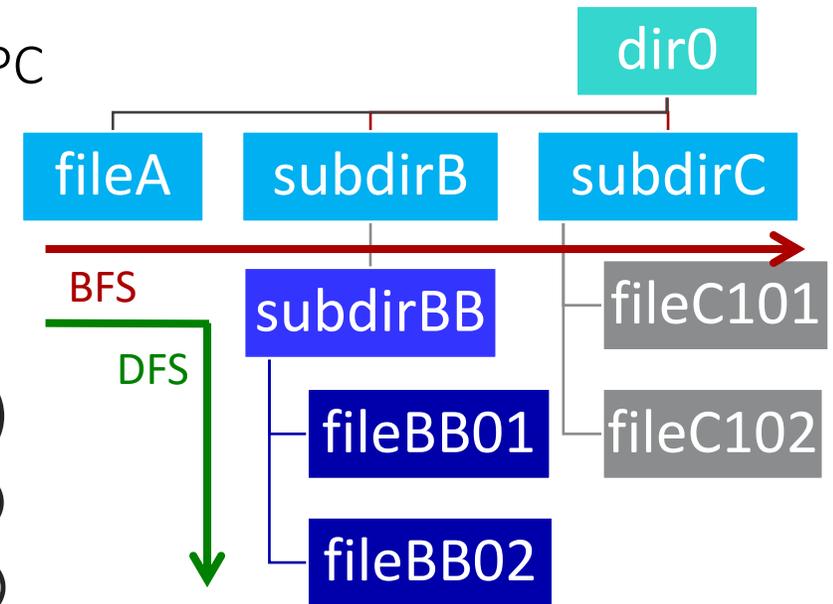
- Allow multiple getattrs/updates packed into a single MDS RPC
- More efficient network and server-side request handling

▶ **Batched statahead** for `ls -l`, `find`, etc. ([LU-14139](#))

- Aggregate getattr RPCs for existing statahead mechanism

▶ **Cross-Directory statahead** pattern matching ([LU-14380](#))

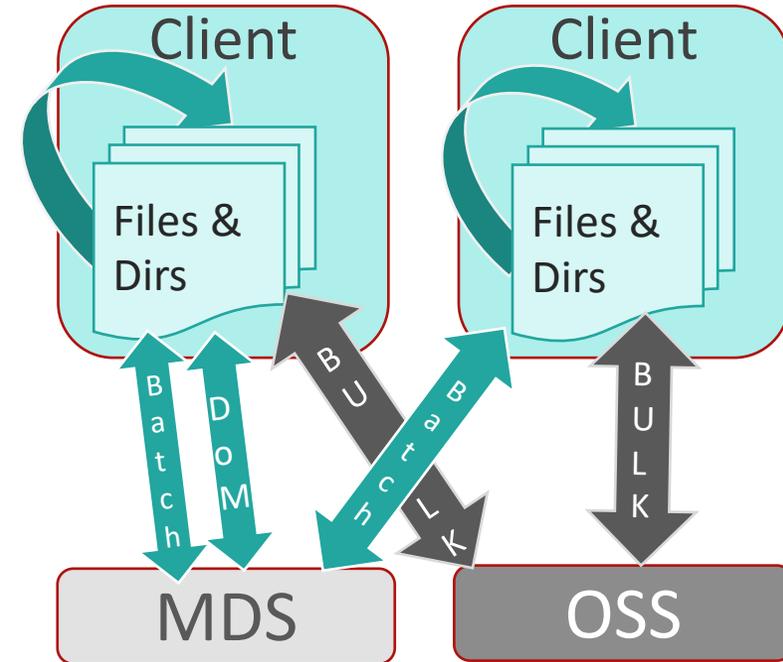
- Existing statahead only detects `readdir()`-ordered `stat()`
- Detect pattern for alphanumeric ordered traversal + `stat()`
- Detect breadth-first (**BFS**) depth-first (**DFS**) directory traversal
- Direct statahead to next file/subdirectory based on pattern



Metadata Writeback Cache (WBC) ([LU-10983](#))

(WC 2.16+) 
Whamcloud

- ▶ Create new dirs/files **in client RAM without RPCs**
 - Lock new directory exclusively at `mkdir` time
 - Cache new files/dirs/data in RAM until cache flush or remote access
- ▶ **No RPC round-trips** for file modifications in new directory
- ▶ Batch RPC for efficient directory fetch and cache flush
- ▶ **Files globally visible on remote client access**
 - Flush top-level entries, exclusively lock new subdirs, unlock parent
 - Repeat as needed for subdirectories being accessed remotely
 - Flush rest of tree in background to MDS/OSS by age or size limits
- ▶ Productization of WBC code well underway
 - Some complexity handling partially-cached directories
 - Need to integrate space usage with quota/grant



MDT DNE Improvements

(WC 2.15+)



▶ **DNE MDT Space Balance** - load balancing with normal `mkdir` ([LU-13439](#), [LU-13440](#))

- Round-robin/balanced subdirs, prefer to stay on parent, limited layout inheritance depth
- Keep MDTs within free inodes/space (`mdt.*.mdt_qos_threshold_rr=5%`)

▶ **Single-dir migration** - "`lfs migrate -m -d <dir>`" ([LU-14975](#))

- Move only one directory level, instead recusing down full subdirectory tree

▶ **Balanced migration** - "`lfs migrate -m -1 <dir>`" ([LU-13076](#))

- Auto-select less-full MDTs for each directory, keep inodes local to parent

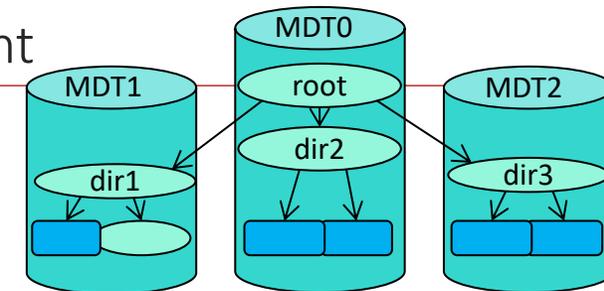
2.15

2.16 ▶ **OST object directory scalability for large OSTs** ([LU-11912](#))

- Request new OST FID Sequences more frequently

▶ **DNE locking, migration, remote RPC optimization** ([LU-15528](#))

- Improve distributed transaction performance, reduce lock contention



Lustre Metadata Redundancy ([LU-12310](#)) (2.17+)

LMR1a: Replicate MDT0000 Services to Other MDS Nodes



▶ Add replication for FLDB across MDTs

- Updated very rarely (new MDT/OST addition, every 1B sequences)
- Maybe OK in LMR1a to only update FLDB when MDT0000 is available?
- Mechanism to mark MDT FLDB as primary copy (version increment?, for LMR1b?)
- Other MDTs fetch primary FLDB and store to local FLDB copy on initial connection

▶ Add replication for Quota Master to other MDTs

- Updated more frequently than FLDB, less replication possible
- Only need to replicate quota *limits*, not usage/accounting (on each target locally)
- Could rebuild quota usage tables on backup MDT from OSTs/MDTs after takeover?

▶ Distribute or failover flock() locking to other MDTs

- Improved performance for independent flock-intensive workloads
- Complex deadlock detection, failover to backup MDT may be easier

LMR1b: DNE Distributed Transaction Performance

- ▶ DNE2 Distributed Transactions have excessive ordering/sync operations
 - Single OSP RPC in flight between each MDT pair (code issue, patch in Gerrit)
- ▶ Compromised from original design due to implementation issues
 - Difficult to journal bit-level updates in recovery llog file
 - Must order updates within a single byte to ensure other bits are not clobbered
- ▶ Use alternative to llog for storing DNE recovery logs
 - Index with named records does not have serialized transaction ordering issues
 - Remove sync write from transaction commit?
 - Batched OSP/OUT updates in a single RPC?
- ▶ Optimizations improve **all** DNE ops, independent of LMR

LMR1c: Replicate Top-level Directories

(2.18+)



- ▶ **ROOT/** directory is replicated across multiple MDTs
 - Up to 7 replicas (?), prefer separate MDTs (different NIDs if possible, also racks?)
 - **ROOT/** is rarely modified, performance overhead should not affect normal usage
 - **ROOT/** is critical for filesystem usage, must always be available to access filesystem
 - **ROOT/** directory lookup by name on secondary MDTs
- ▶ Read-only under downgrade mode (only directory replication, not files?)
- ▶ New incompatible flag `LMAI_REPLICA` in LMA EA
- ▶ User tools to add/remove replicas to directories on non-LMR system
- ▶ Enhance LFSCK to verify/repair consistency between replica objects

- ▶ Additional LMR2/3 phases to reach full MDT redundancy



Whamcloud

Thank You!
Questions?