# Lustre Feature Renaissance

Nathan Rutman
Lug 2018

# 14 Years of Scaling

**CRAY**

GB/s

- NERSC NERSC9
- LANL Trinity (#5)
- NCSA BlueWaters
- Riken K (#1)
- LLNL Sequoia (#1)
- ORNL Spider 2
- MetOffice Cray XC40 (#11)
- Argonne Aurora
- ORNL Titan (#1)
- NUDT Tianhe-2 (#1)
- NERSC Cori (#8)
- 1000
- CEA Tera 100 (#9)
- TOTAL Pangea (#11)
- Kaust Shaheen II (#7)
- Mistral
- PACS Oakforest (#9)
- HLRS Hazel Hen (#8)
- TACC Stampede2 (#12)
- ORNL Jaguar (#1) LLNL aggregate
- Argonne Mira (GPFS) (#3)
- NRCPC Sunway TaihuLight (#1)
- Argonne Theta (#16)
- LANL Cielo (#6)
- TACC Stampede (#6)
- TACC Stampede (#6)
- CSCS Piz Daint (#6)
- 100
- CEA Tera 10 (#7)
- NOAA Gaea (#32)
- NUDT Tianhe-1A (#1)
- TiTech Tsubame 1 (#7)
- NASA Pleiades (#7)
- Sandia Red Storm (#6)
- ORNL Kraken (#8)
- TACC Ranger (#2)
- LLNL BlueGene/L (#1)
- Sandia Red Sky (#10)
- NERSC/LBNL Franklin (#9)
- Indiana Big Red (#42)
- ...ngsten (#20)
- 10
- CINECA Marconi (GPFS) (#12)
- NCSA Abe (#14)
- LLNL Thunder (#2)
- MCR (#3)
- 1
- TACC Lonestar (#12)
- 2005
- 2010
- 2015
- 2020

# Years of effort spent on performance and scaling

Not an exhaustive list
- ldiskfs scaling, ZFS
- Recovery: interop, VBR, AT, FSCK, Imperative Recovery
- CLIO, MDS rewrite, FIDs
- IO scaling: LRU, read cache, readahead, wide striping, multi-MB rpc, DoM
- MD scaling: statahead, DNE, MMR
- tons of diverse performance improvements
- bugs bugs bugs

Features too of course: mountconf, Kerberos, NRS, HSM, changelogs, pools, multirail

# Capable infrastructures in place

- DNE – MD horizontal scaling
- Complex layouts – much more interesting data placement
- FLR – data redundancy inside Lustre

## Time to reap some Feature rewards

- Let's look at some possible features
- These aren't even designs, just ideas
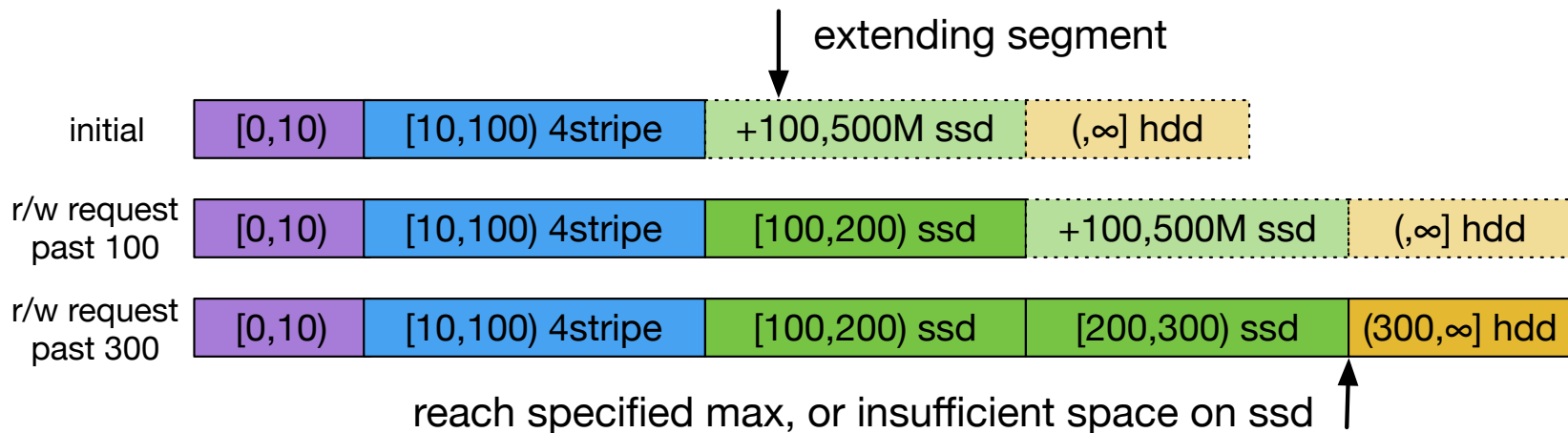
# FLR – data redundancy *inside* Lustre

- **One small (?) step for Layouts, one giant leap for Lustre systems design**
- **No longer need to rely on ~~Failover~~ for data access**
- **Dual-ported, dual-server, dual-path, dual-$ - nope.**
- **Need:**
  - FLR2 = immediate – client writes data durably
  - FLR3 = EC – boo 100% overhead, yay 20% overhead
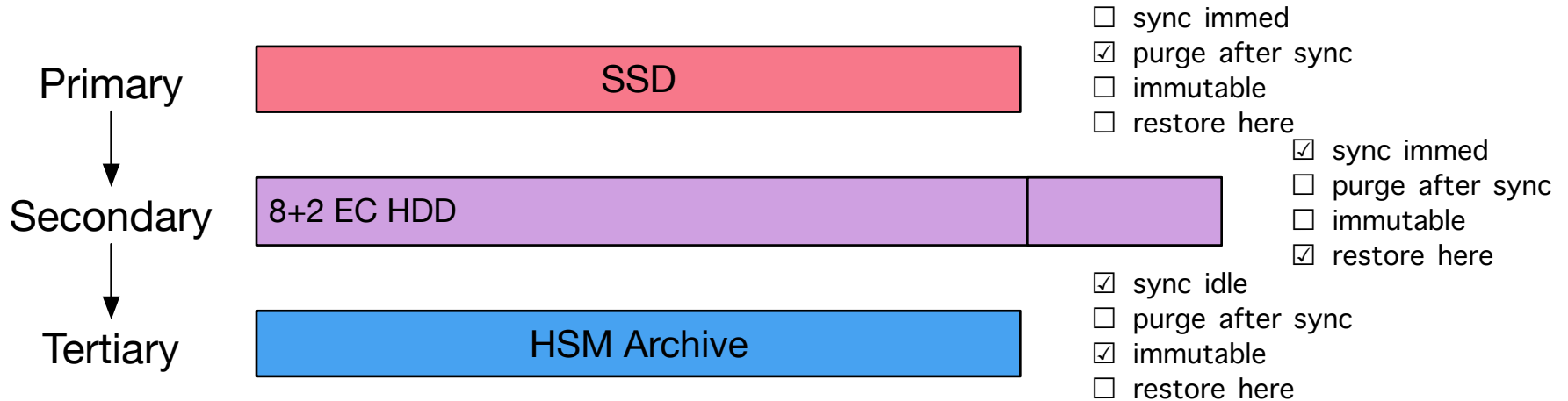  - Degraded write support. Track changes for reintegration, or asymmetric layouts?

# Spillover Space: death to ENOSPC

- **Self-extending PFL** (LU-10070, LU-10169)
- **Some PFL segments are virtual, instantiated on demand**
- **Request in a virtual segment requires layout update**
- **MDS adds a new component on demand**
- **Can choose the new component striping based on dynamic conditions (e.g. free space)**

extending segment

| initial | [0,10) | [10,100) 4stripe | +100,500M ssd | (,∞] hdd |

| r/w request past 100 | [0,10) | [10,100) 4stripe | [100,200) ssd | +100,500M ssd | (,∞] hdd |

| r/w request past 300 | [0,10) | [10,100) 4stripe | [100,200) ssd | [200,300) ssd | (300,∞] hdd |

reach specified max, or insufficient space on ssd

# ILM Layouts

- **Layout implies an action:** stale FLR copy = resync w/ lfs mirror
- **And a timeframe: (immediate | eventual)**
- **Simple ILM policy** *already* **encoded into layout**
- **Add some flags to layout and/or policy ref**
- **Make HSM a true layout** (LU-10606): stale HSM copy = resync w/ lfs hsm
- **Use Coordinator and Copytool for all movement** (LU-6081)



Primary → Secondary → Tertiary

SSD
- ☐ sync immed
- ☑ purge after sync
- ☐ immutable
- ☐ restore here

8+2 EC HDD
- ☑ sync immed
- ☐ purge after sync
- ☐ immutable
- ☑ restore here

HSM Archive
- ☑ sync idle
- ☐ purge after sync
- ☑ immutable
- ☐ restore here

# Asymmetric Layouts

- **Reads go to R iff not in W**
- **Block bitmap on W tracks newly written data**
- **Client caching and DIO insure full-page writes**
- **W controls all locks, gives bitmap along with lock grants**
- **Clients access R or W directly, all under W's locks**

**Why?**
- Write to flash, read from HDD
- Continue writing to new W if an OST fails (checkpoint) (or ENOSPC)
- EC degraded write case – point of EC is to remain usable in failures

Asym Layout

write

read

Block bitmap

W layer

R layer

# Fast Find

- **Why do we copy Lustre MD into DB's or scan raw ldiskfs?**
- **Need to quickly find files that match certain criteria**
- **A great 'lfs find' could do the same thing, saving the tools effort**
  - Server side. RPC from client, returns filtered list
  - Logical combinations of filters
  - Unix-style piping:  *lfs find /lustre -size +20M | lfs hsm archive*
- **Add new MDT indices to efficiently generate initial candidate lists**
  - LRU, file heat, mtime, size
  - dt_index_operations (eg IAM) provides generic indexing code
  - Update indices transactionally with MD updates

# Rough SoM

- **FLR records file size on MDS; comes with sync**
- **DoM records file size on MDS**
- **Straightforward to get maximum size, if we don't care about evicted/failover case**
- **Rough size is fine for many purposes (e.g. policies)**
- **Record the quality of SoM, let users decide if usable**
- **Strict, Rough, Stale, Unknown (LU-9538)**
- **Don't return as POSIX size unless strict**
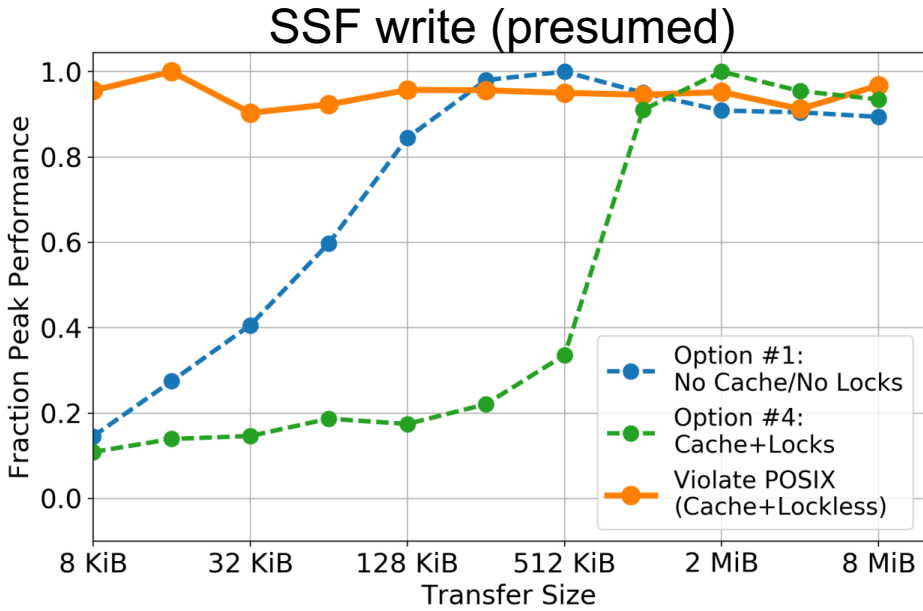
# Clone Files: extreme create scaling



- **File create: ask MDS to create, lock dir, create inode, assign objects**
- **Clone create: create "all" the files at once**

- **Single MDS inode, single namespace entry: foo.#**
- **FID is prefix+#**
- **Layout is f(FID)**

| MDS Inode /pdir/foo.# |
|---|
| uid,gid |
| perms |
| pdir |
| FID ABCD00xxxxx |

| Layout |
|---|
| canonical_ost_list |
| stripe_count, size |
| selected_osts = f(FID, c_o_l) |

| Client Inode /pdir/foo.12345 |
|---|
| uid,gid |
| perms |
| pdir |
| FID ABCD0012345 |

| Layout |
|---|
| canonical_ost_list |
| stripe_count, size |
| selected_osts = f(FID, c_o_l) = {x, y} |

OST x

OST y

- **Shared MD (clones!) but different objects / data / sizes**
- **open(foo.4,O_CREAT) is now a client-local operation**

# Alternate Consistency Models

- **POSIX API vs POSIX consistency semantics**
- **Caching allow write coalesce, local latency**
- **But pay a penalty for locks**
- **Solutions in Lustre, but requires effort**
  - Lockless DIO, Grouplocks
- **Make it easier**
  - Iadvise?
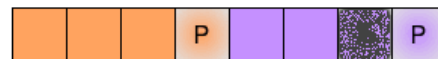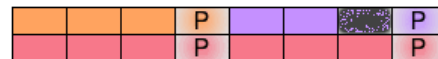  - Persistent file tags?
  - Automatically change modes?

### SSF write (presumed)



Glenn Lockwood
https://www.nextplatform.com/2017/09/11/whats-bad-posix-io/

# All Together Now

# Feature Vote?

1. FLR EC with degraded writes
2. ~~Spillover Space~~ *Cray*
3. ILM Layouts
4. Asymmetric Layouts
5. Fast Find
6. ~~Rough SoM~~ *DDN*
7. Clone Files
8. Alternate Consistency Hints

# Implementation Plan

1. **Ignore Nathan's slideware**
2. **<insert smart developer here>**



3. **Implement!**