A quantitative approach to architecting all-flash Lustre file systems





Glenn K. Lockwood et al.

June 20, 2019





NERSC's mission and workload mix

- NERSC is the mission HPC computing center for the DOE Office of Science
 - HPC and data systems for the broad Office of Science community
 - 7,000 Users, 870 Projects, 700 Codes
 - >2,000 publications per year
- 2015 Nobel prize in physics supported by • **NERSC** systems and data archive
- Diverse workload type and size:
 - Biology, Environment, Materials, Chemistry, Geophysics, Nuclear Physics, Fusion Energy, Plasma Physics, Computing Research
 - New experimental and AI-driven workloads

- 2 -









NERSC's 2020 System: Perlmutter





- Designed for large-scale simulation and large-scale data analysis
- 3x-4x capability of current system
- Include both NVIDIA GPUaccelerated and AMD CPU-only nodes
- 200 Gb/s Cray Slingshot interconnect
- Single tier, all-flash Lustre file system







- 4 -

Designing an all-flash parallel file system

- In 2020, all-flash makes good sense for the performance tier
- With a <u>fixed budget</u>, how much should we spend on...
 - OST capacity?

Office of

- NVMe endurance?
- MDT capacity for inodes and Lustre DOM?
- Performance not discussed here





NERSC's quantitative approach to design



• Inputs:

- 1. Workload data from a reference system
- 2. Performance specs of future compute subsystem
- 3. Operational policies of future system
- Outputs:
 - Range of minimum required OST capacity
 - Range of minimum required SSD endurance
 - Range of minimum required MDT capacity for DOM





Reference system: NERSC Cori

Compute

- 9,688 Intel KNL nodes
- 2,388 Intel Haswell nodes

Storage

- 30 PB, 700 GB/s scratch
 - Lustre (Cray ClusterStor)
 - 248 OSSes x 41 HDDs x 4 TB
 - 8+2 RAID6 declustered parity

• 1.8 PB, 1.5 TB/s burst buffer

- Cray DataWarp
- 288 BBNs x4 SSDs x 1.6 TB
- RAIDO







NERSC workload data sources





Science

- smartmon tools
 - Per-device metrics vs time
 - Write amplification factor
 - Total bytes written

Lustre Monitoring Tools (LMT)

- Per-OST measurements vs time
- Bytes written to OSTs
- Lustre "Ifs df" & cron
 - Per-OST fullness vs time
 - LMT also has this information

Robinhood

- Snapshot of entire namespace & POSIX metadata
- File size distribution
- Non-file inode sizes



Minimum required file system capacity

















How much flash capacity do we need?



- 1. How quickly **users fill** the file system
- How frequently the facility drains the file system









File system capacity model





Calculating *C*^{new} for Perlmutter





Mean daily growth for Cori is <u>133 TB/day</u>



Data retention policy for Perlmutter:

- Purge/migrate after <u>28 days</u>
- Remove/migrate <u>50% of</u> <u>total capacity</u>



SSI Perlmutter will be <u>3x to 4x</u> "faster" than reference system

Office of

Science

C^{new} Minimum capacity is between <u>22 PB and 30 PB</u>



Do we need [to pay for] highendurance SSDs?



















Intuitively, required drive writes per day depends on...

- 1. How much data users write to the file system daily
- 2. How much extra RAID parity is written
- 3. How much hidden "stuff" is written
 - 1. Read-modify-write write amplification
 - 2. Superblock updates
- 4. How big a "drive write" actually is





Calculating drive endurance requirements





Science



ERKELEY LA

Cori's FSWPD over two years

From LMT/pytokio

- 1 FSWPD = 30.5 PB
- Mean FSWPD 0.024 • (~700 TB/day)

Office of ENERG Science

U.S. DEPARTMENT OF Office of Science

Estimating WAF in production

- <u>Not</u> internal SSD WAF
- Difference between what Lustre writes and what disks write
 - Can use OST-level writes (LMT) and SMART data
 - Had to use DataWarp
 SSD WAFs here
- Median WAF: <u>2.68</u>
- 95th percentile: <u>3.17</u>











How much MDT capacity is needed for DOM?















directly on Lustre MDT

- Reduces number of RPCs, small-file interference, etc
- But HPC facilities now have to define:
 - How big should S₀ be?
 - How much MDT capacity is required to store all files' DOM components?
 - How much MDT capacity is required to store inode data?













Intuitively, MDT capacity depends on...

- 1. How much capacity we need to store inodes (files, directories, symlinks, ...)
- 2. How many files we have with size smaller than S₀
- 3. How many files we have with size larger than S₀
- 4. Our choice of S_0



Office of Science

- 22 -

Fraction of non-file inodes

10

MDT capacity required for inodes

- 1 inode = 4 KiB
- ...but big dirs need more MDT capacity!
- Assume size dist is constant for any C^{new}
- Then estimate the size dist for our C^{new}





MDT capacity required for DOM









Total MDT capacity required vs. S₀

- Uncertainty from scaling distributions between Cori and Perlmutter
- Extreme values of S₀ uninteresting
 - Tiny S₀ = nothing fits in DOM
 - Huge S₀ = everything fits in MDT (and I/O is no longer parallel)







Cost-performance tradeoff of DOM









Conclusions



We have defined models to quantify relationships between

- Workload (I/O and growth rates, file size distributions)
- Policies (purge/data migration policy)
- File system design parameters
 - File system capacity
 - SSD endurance
 - MDT capacity and DOM configuration

Imperfect models, but still

- identify key workload parameters
- bound on design requirements based on facts, not instincts
- serve as a starting point for sensible system design

All code and workload data available online!*

https://doi.org/10.5281/zenodo.3244453







Thank you!

(and we're hiring!)





