

# Profiling Platform Storage Using IO500 and Mistral

Nolan Monnier, **Jay Lofstead**,  
Margaret Lawson, Matthew Curry



*Exceptional  
service  
in the  
national  
interest*

PDSW 2019

SAND2019-14085 C



Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-NA0003525.

# Astra Testbed

- Prove viability of advanced technologies for NNSA integrated codes, at scale
- Expand the HPC-ecosystem by developing emerging yet-to-be proven technologies
- Is technology viable for future ATS/CTS platforms supporting ASC mission?
- Increase technology AND integrator choices
- Buy down risk and increase technology and vendor choices for future NNSA production platforms
- Ability to accept higher risk allows for more/faster technology advancement
- Lowers/eliminates mission risk and significantly reduces investment
- Jointly address hardware and software technologies
- First Prototype platform targeting Arm Architecture
- **#156 on top500 (June 2019) (2.2 PF)**

# Astra Hardware

**HPE Apollo 70 Chassis: 4 nodes**



**HPE Apollo 70 Rack**



**18 chassis/rack**

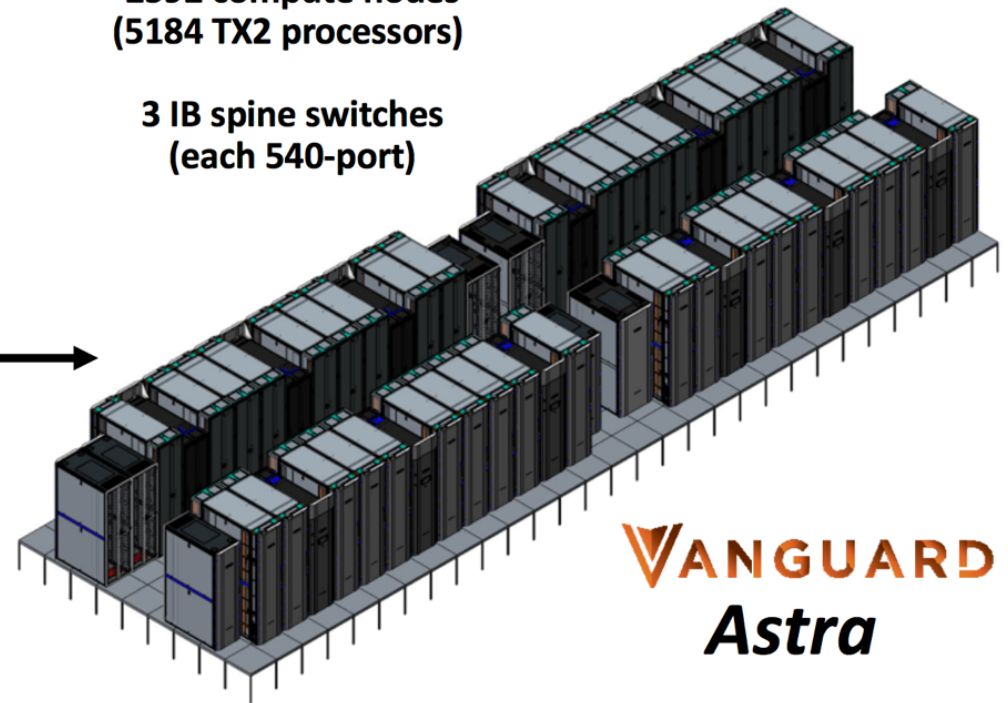
**72 nodes/rack**

**3 IB switches/rack**  
(one 36-port switch  
per 6 chassis)

**36 compute racks**  
(9 scalable units, each 4 racks)

**2592 compute nodes**  
(5184 TX2 processors)

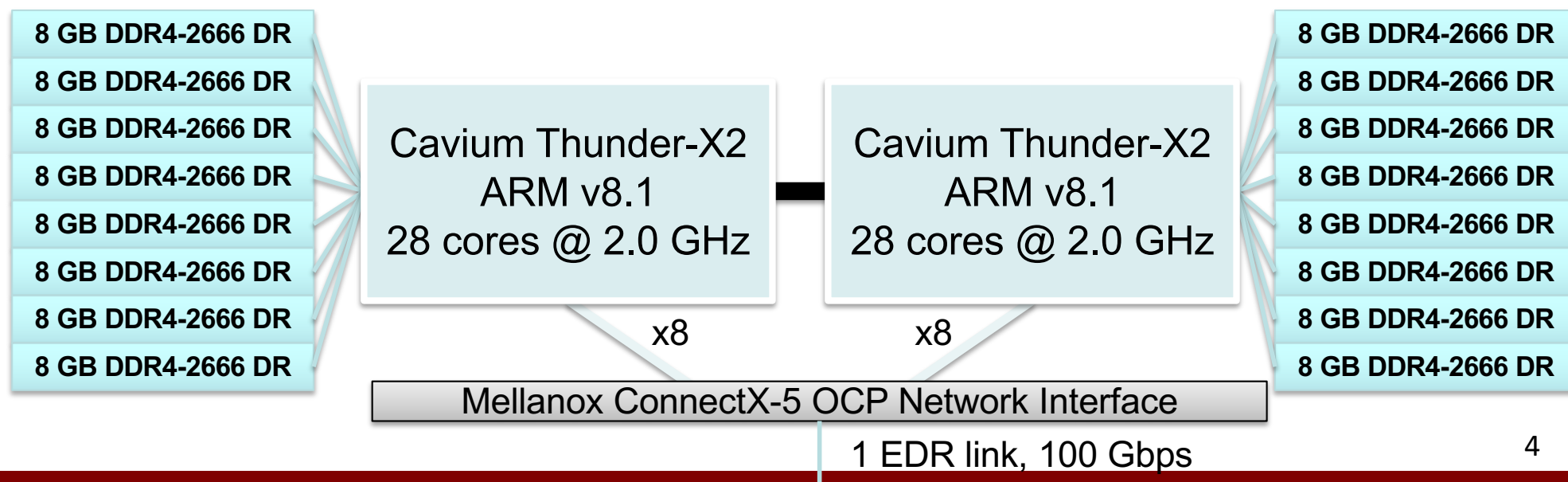
**3 IB spine switches**  
(each 540-port)



**VANGUARD**  
*Astra*

# Astra System Specs

- 2,592 HPE Apollo 70 compute nodes
- Cavium Thunder-X2 Arm SoC, 28 core, 2.0 GHz
- 5,184 CPUs, 145,152 cores, 2.3 PFLOPs system peak
- 128GB DDR Memory per node (8 memory channels per socket)
- Aggregate capacity: 332 TB, Aggregate Bandwidth: 885 TB/s
- Mellanox IB EDR, ConnectX-5
- HPE Apollo 4520 All-flash storage, Lustre parallel file-system
- Capacity: 403 TB (usable)
- Bandwidth 244 GB/s



# Astra Storage Specs

- `/lustre` filesystem
  - 40 OSSes
  - 2 MDSes
  - 21 1.6TB NVMe devices per OSS spread across 3 ZFS pools per node using raidz
  - 240 GB/sec peak bandwidth
  - 990 TB usable storage
- `/oscratch` filesystem
  - 8 OSSes, each serving 10 OSTs
  - 2 MDSes, 1 MDT each
- Storage has full bandwidth to interconnect core (4xEDR 100 Gpbs Infiniband)

# IO500

- Virtual Institute for IO (<https://www.vi4io.org>) hosted
- <http://io500.org>
- Collect details information about storage and platforms
- Provide a balanced way to compare storage systems
  - No, it is not perfect
- MDTest, IOR, and find
- “hard” tests for worst case scenarios
- “easy” tests user configurable to showcase a system’s potential
- Find to represent walking the file tree for purge or similar
- Geometric mean of all values to get resulting score
- Published at ISC and SC every year (Tuesday, 12:15, 205-207)

# Profiling IO

- Darshan is available, but we wanted to see how this Ellexus' Mistral tool worked and then compare
  - Future work
- Mistral offers node-level statistics
  - Breeze is the per process tool
- Collect on a per-second basis
- Ideal configuration is Elasticsearch and Grafana
  
- How well is our ARM-based Lustre client and the back end Lustre system working?
- Can we learn using IO500's setup instead of apps or app proxies?

# Tuning IO500 to Push the System

- 1) Obtain system information and theoretical characteristics.
- 2) Set test directories' stripe size based on test files' size and number of storage targets.
- 3) Determine number of nodes to use.
- 4) Increase the cores per node to maximize bandwidth, until the bandwidth for ior easy reasonably approaches a theoretical limit.
- 5) Adjust the cores per node to balance bandwidth and metadata results



# Tuning Final Configuration

- Configured striping to best suit the "hard" tests
- Set nodes used == to number of storage targets
  - `/lustre` used 58 procs (on the 58 cores) on 121 nodes
- 10 node challenge ran proportional configurations
- Metadata did not require special stripe configurations

# IO500 Scores

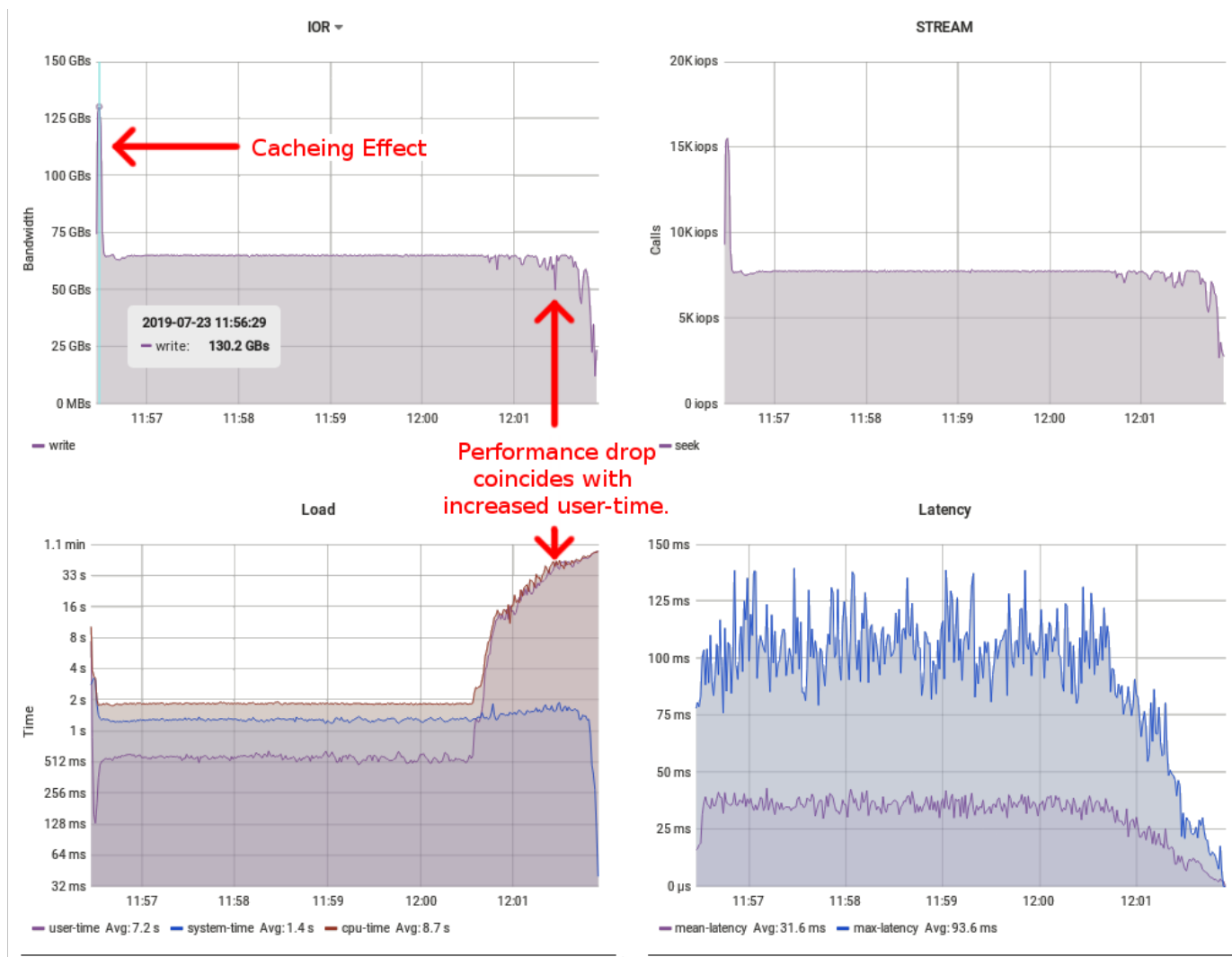
Nodes	Mount Point	Bandwidth	IOPS	Score
121	/lustre	84.8118 GB/s	35.9847 kiops	55.2443
10	/lustre	28.4097 GB/s	45.7227 kiops	36.0412

- 14<sup>th</sup> on the ISC19 overall list
- 13<sup>th</sup> on the ISC19 10 node challenge
- Plug: (see new results at IO500 BoF Tuesday 12:15, 205-207)

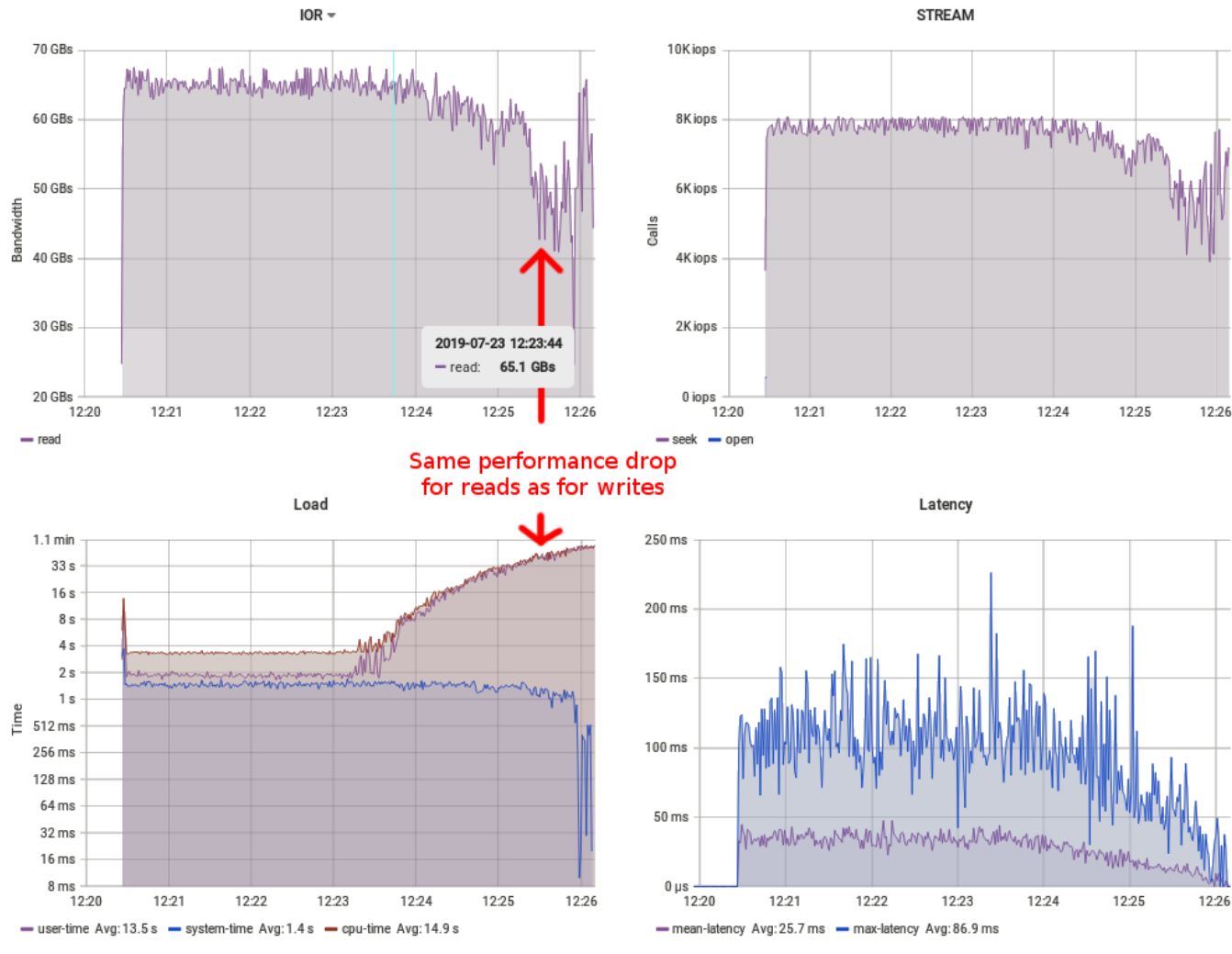
# Problems!

- On `/oscratch` filesystem, we ran out of space in directory when creating files
  - Lustre has an option to extend capacity, but we were not allowed to change it.
- This prevented getting “official” `/oscratch` results

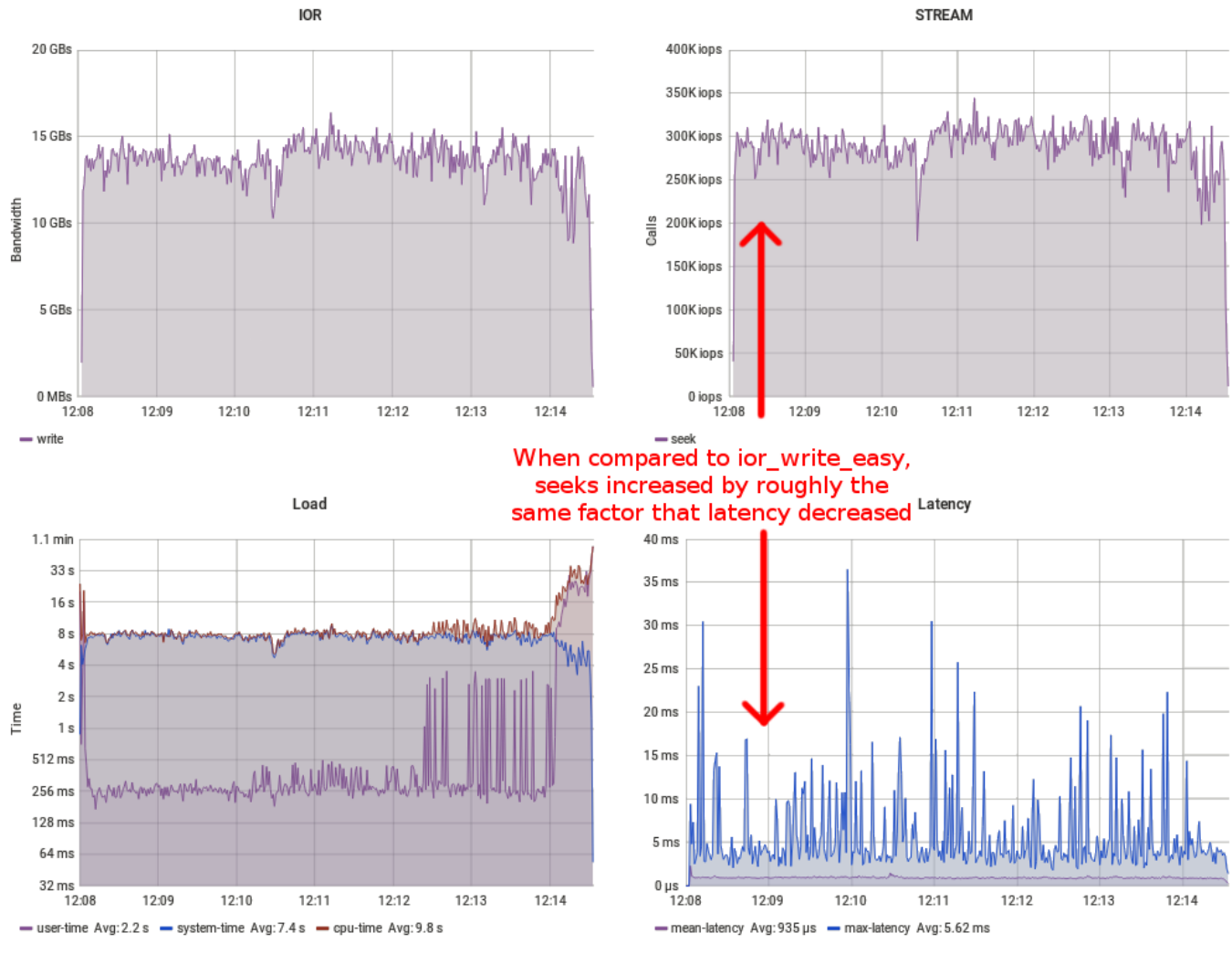
# IOR Write Easy



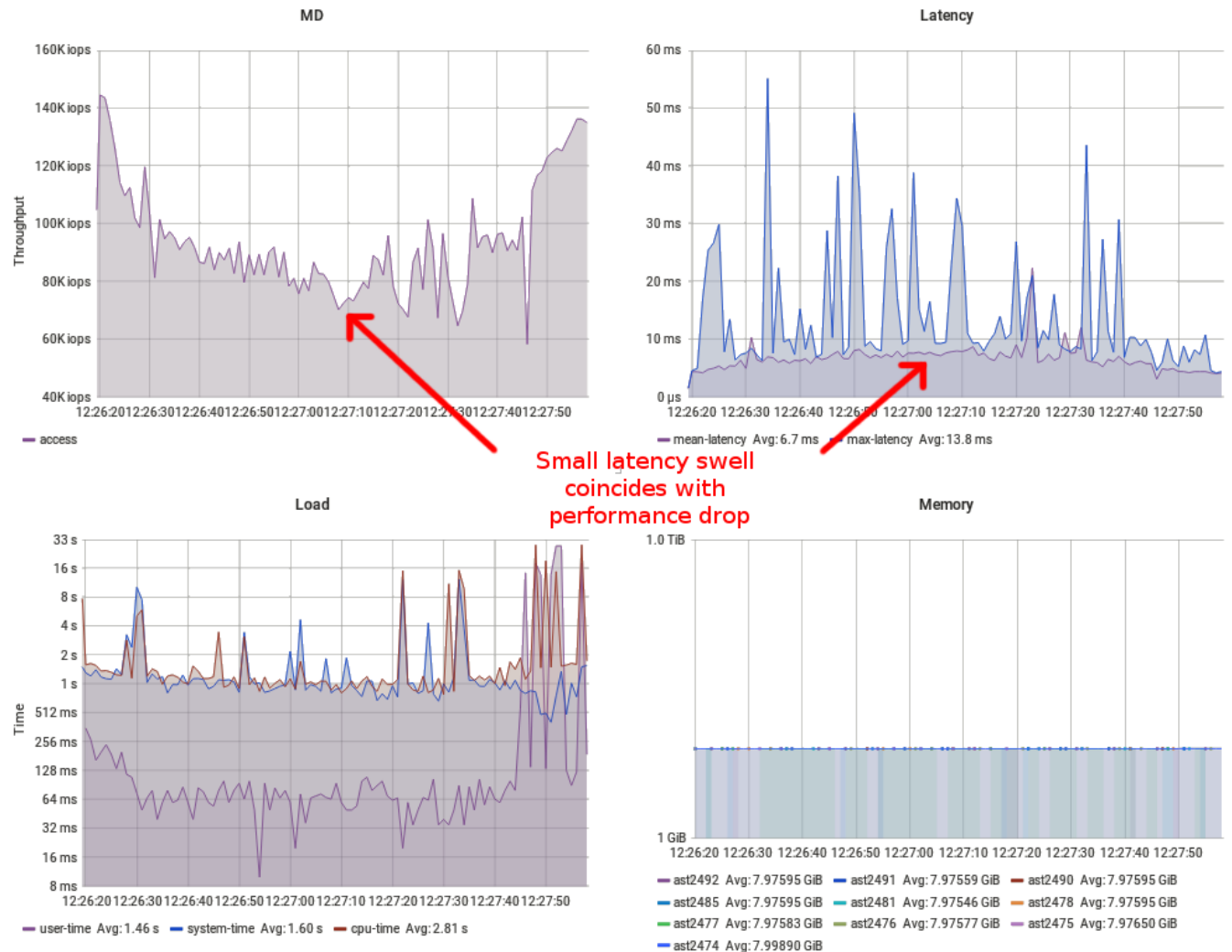
# IOR Read Easy



# IOR Write Hard



# MD Stat Easy

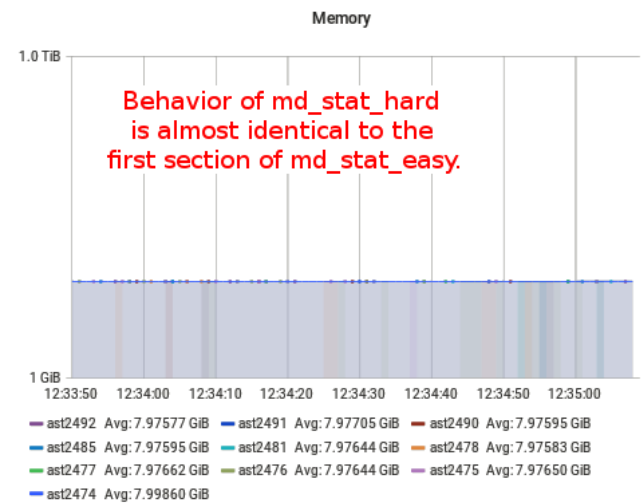
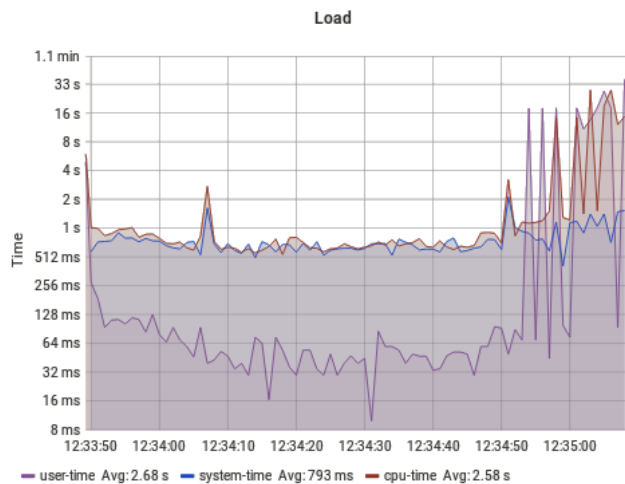
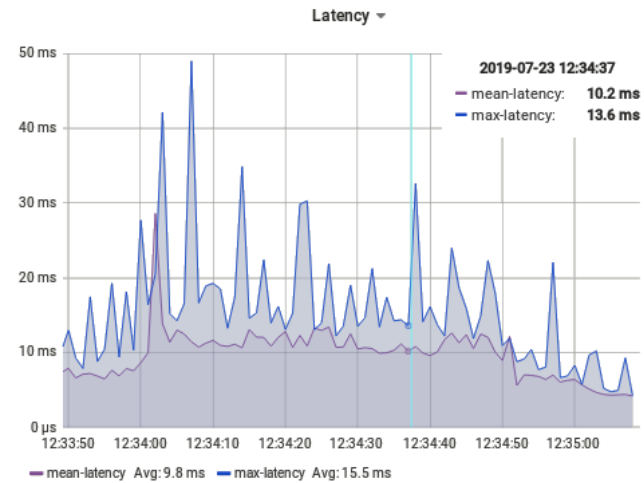
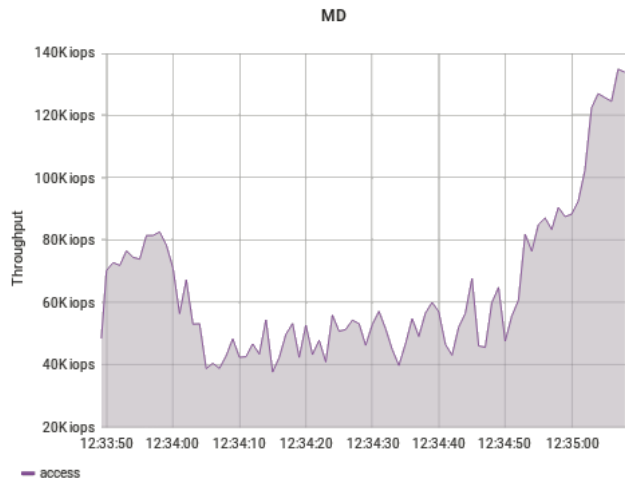


# MD Stat Easy

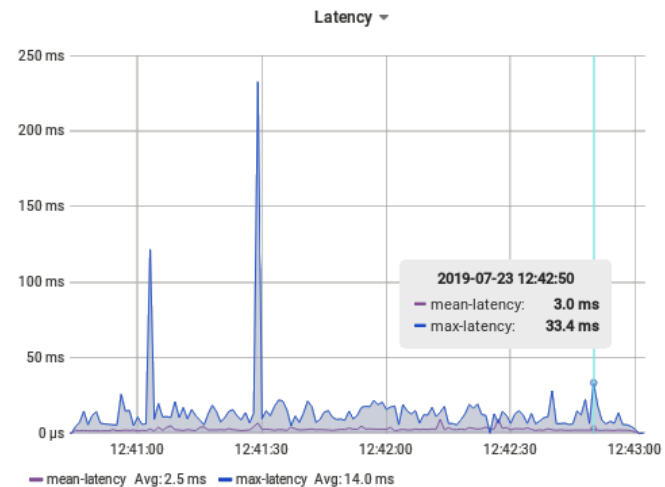
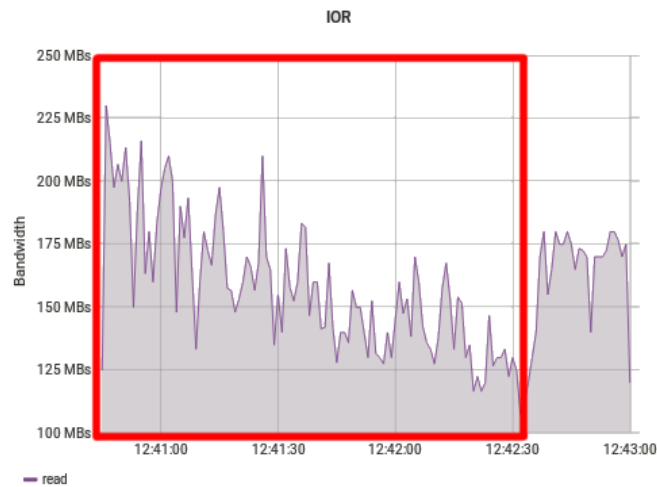
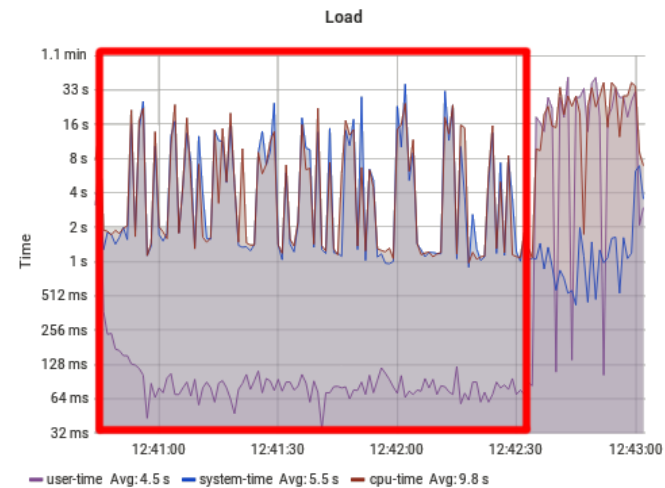
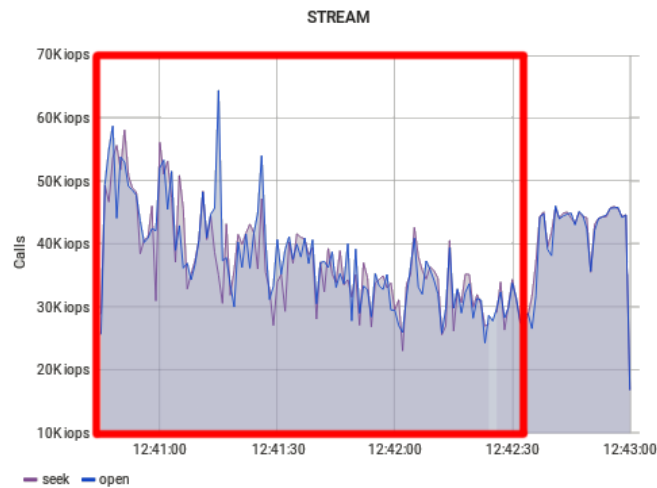




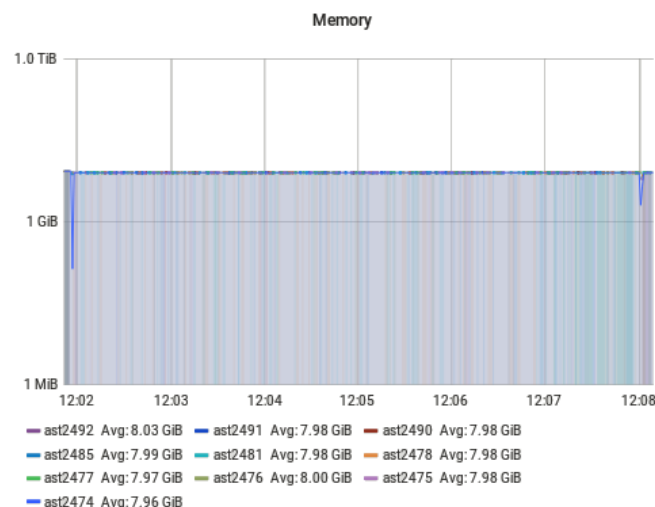
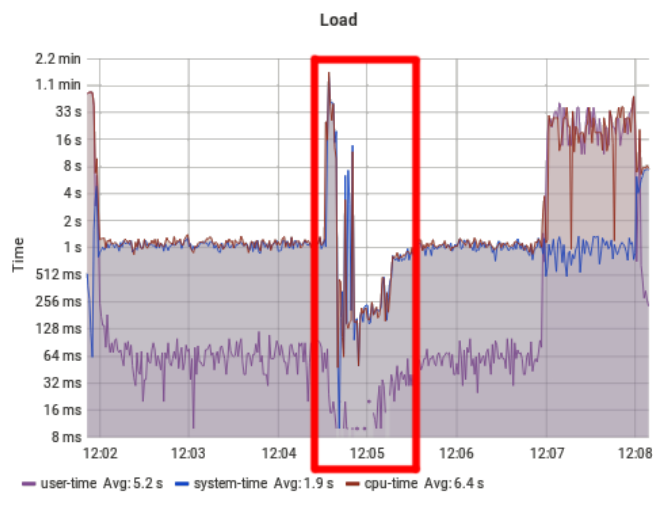
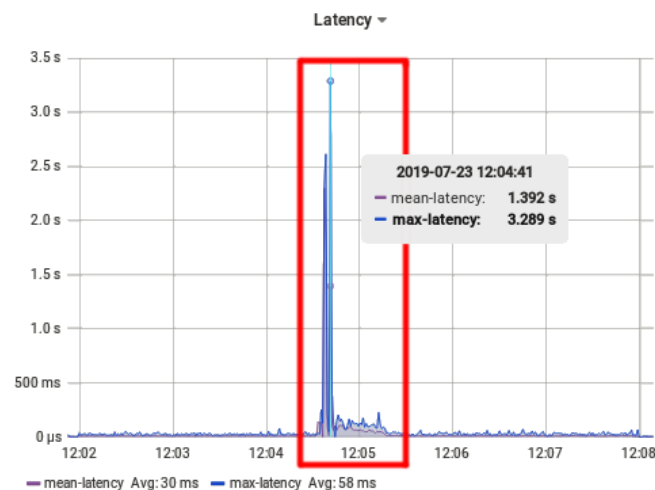
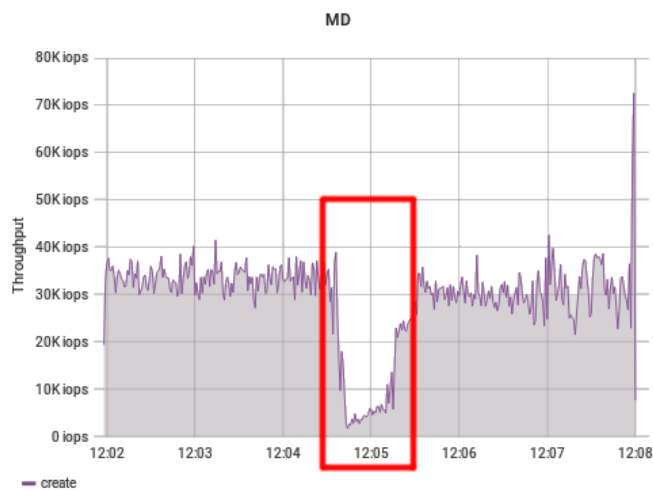
# MD Stat Hard



# MD Read Hard



# MD Write Easy



# Conclusion

- Hit 33% of peak (midpoint of 20%-50% of peak expected)
- Reproducibly issues with the Lustre client and MPI we had not isolated before
- Found gradual performance degradations and odd performance fluxuations to debug

# Acknowledgements

- Thank you to Ellexus for supporting this work
- Supported under the Vanguard project from NNSA
- Sandia National Laboratories is a multimission laboratory managed and operated by National Technology and Engineering Solutions of Sandia, LLC, a wholly owned subsidiary of Honeywell International, Inc., for the U.S. Department of Energys National Nuclear Security Administration under contract DE-NA0003525.