# outthink limits

# Spectrum Scale Enhancements for CORAL

**Sarp Oral, Oak Ridge National Laboratory**
**Gautam Shah, IBM**

SC16
Salt Lake City. | hpc
Utah | matters.

# What is CORAL

Collaboration of DOE Oak Ridge, Argonne, and Lawrence Livermore National Labs

- Established in early 2014 to leverage supercomputing investments, streamline procurement processes and reduce costs to develop supercomputers
  - "High-performance computing is an essential component of the science and technology portfolio required to maintain U.S. competitiveness and ensure our economic and national security" - U.S. Secretary of Energy Ernest Moniz
- Two new High Performance Computing (HPC) awards announced in November 2014
  - Both CORAL awards leverage the IBM Power Architecture, NVIDIA's Volta GPU and Mellanox's Interconnected technologies to advance key research initiatives for national nuclear deterrence, technology advancement and scientific discovery
    - Oak Ridge National Laboratory's (ORNL's) new system, Summit, is expected to provide at least five times the performance of ORNL's current leadership system, Titan
    - Lawrence Livermore National Laboratory's (LLNL's) new supercomputer, Sierra, is expected to be at least seven times more powerful than LLNL's current machine, Sequoia.

Source: http://energy.gov/articles/department-energy-awards-425-million-next-generation-supercomputing-technologies

# CORAL Systems

- **LLNL's Sierra system**
  - ~4000 Power9 nodes with GPU acceleration
  - ~2.3 PB system memory (include DDR & HBM; does not include NVMe)
  - Dual-rail InfiniBand EDR fat tree network or better
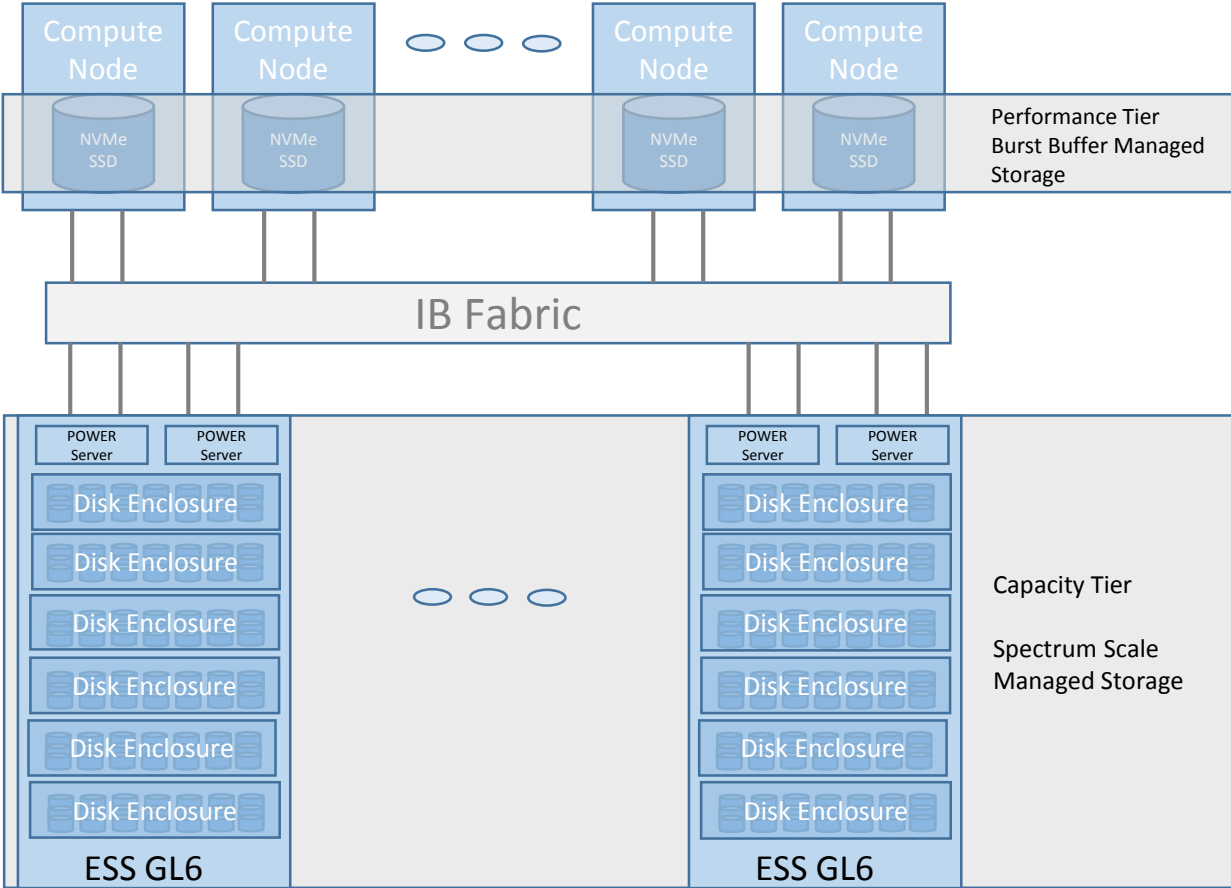  - ~120 PF
  - ~9 MW
- **ORNL's Summit system**
  - ~4500 Power9 nodes with GPU acceleration
  - ~ 2.7 PB system memory (include DDR & HBM; does not include NVMe)
  - Dual-rail InfiniBand EDR fat tree network or better
  - ~200 PF peak
  - ~13 MW

# CORAL System Storage Overview

Storage architecture

- Need for a burst buffer/performance tier
    - Lowers traditional spinning disk & lower power consumption
    - Node local NVMe SSD managed by Burst Buffer Software
- Capacity requirement – ESS Storage
- Performance/Scaling requirements – Spectrum Scale Software

# CORAL Storage Overview

# Burst Buffer Software

- Goals/Features
  - Support node-local checkpoints
    - SSD partitioned and formatted with standard Linux file system
  - Support staging data in and out of SSD
    - Provide a mechanism for pre- and post- job transfers for data staging via LSF
  - Build asynchronous file transfer service between SSD and GPFS
    - Software on compute node initiates transfer and can poll for completion
  - Avoid excessive data movement
  - Avoid performance jitter to running applications on compute node
    - Move data between compute node and ESS with NVMe over Fabrics for low performance impact
  - SSD wear awareness, health monitoring, and protection

# Spider 3 @ OLCF

Spider 3 is a center-wide single namespace POSIX file system to serve all OLCF resources eliminating data islands, and enabling seamless data sharing between resources

- Built on IBM's Elastic Storage Server based on Power 9 Processor and uses Spectrum Scale (formerly known as GPFS) parallel filesystem technology utilizing GPFS Native RAID with 8+2 redundancy
- Provides a usable capacity of 250 PB
- Performs at an aggregate sequential peak read/write bandwidth of 2.5 TB/s
- Performs at an aggregate random peak read/write bandwidth of 2.2 TB/s
- Provides rich metadata performance; single directory parallel create rate of 50,000/s
- Provides rich interactive performance; @32 KiB I/O 2.6 million IOPs
- Disk-based, with tens of thousands of disks
- Connected to OLCF's SION 3 SAN with IB EDR
- Will also serve as the Summit Burst Buffer sink and source on the end-to-end I/O path

# Spectrum Scale Enhancements for Scaling Namespace

The single namespace CFS will meet the following

- Single name space supporting 250 PB capacity
- Total number of files supported is 100 B
- Maximum file size equal to aggregate system memory
- 10 M files per directory

Enhancements needed in Spectrum Scale

- Improvements in fsck – time to run (including nodes to use), progress reporting, …
- Parallel virtual disk creation
- Reduce contention to allow more concurrency

# Spectrum Scale Enhancements for Scaling Performance

Performance improvement are required to meet:

- Aggregate sequential peak read/write bandwidth of 2.5 TB/s
- Aggregate random peak read/write bandwidth of 2.2 TB/s
- Single directory parallel create rate of 50,000/s
- Interactive performance; @32 KiB I/O 2.6 million IOPs

Enhancements needed in Spectrum Scale:

- Performance counters to help uncover bottlenecks (mmfsadm dump iocounters/iocountercpu)
- Improve RPC communication (avoid global receive pool mutex by creating multiple pools of worker threads;  Fast Condition Variable for some of the condition variables in RPC path; spread interrupt load across multitple IRQs using RDMA completion vectors)
- Improve parallelism of full track writes

# Spectrum Scale Enhancements for Scaling Metadata Rate

Metadata requirements include:

• Single directory parallel create rate of 50,000/s

• Interactive performance; @32 KiB I/O 2.6 million IOPs

Enhancements needed in Spectrum Scale:

• Improve directory block management (avoid directory fragments, directory block split option,

• Avoid token manager revokes

• Pre-allocation of directory blocks (mmchattr)

• Smooth the filesystem sync work over the sync period

# Collaboration for Success

We expect other challenges we have to overcome as we deliver/deploy the system and we are working together to anticipate and resolve these issues

- Impact of rebuild performance over the population size …
- Identify and eliminate "slow" disks so the system performance consistently

# Thank you!



**ibm**.com/systems/hpc

# Legal notices

# Information and trademarks

# Special notices

This document was developed for IBM offerings in the United States as of the date of publication. IBM may not make these offerings available in other countries, and the information is subject to change without notice. Consult your local IBM business contact for information on the IBM offerings available in your area.

Information in this document concerning non-IBM products was obtained from the suppliers of these products or other public sources. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products.

IBM may have patents or pending patent applications covering subject matter in this document. The furnishing of this document does not give you any license to these patents. Send license inquires, in writing, to IBM Director of Licensing, IBM Corporation, New Castle Drive, Armonk, NY 10504-1785 USA.

All statements regarding IBM future direction and intent are subject to change or withdrawal without notice, and represent goals and objectives only.

The information contained in this document has not been submitted to any formal IBM test and is provided "AS IS" with no warranties or guarantees either expressed or implied.

All examples cited or described in this document are presented as illustrations of  the manner in which some IBM products can be used and the results that may be achieved. Actual environmental costs and performance characteristics will vary depending on individual client configurations and conditions.

IBM Global Financing offerings are provided through IBM Credit Corporation in the United States and other IBM subsidiaries and divisions worldwide to qualified commercial and government clients. Rates are based on a client's credit rating, financing terms, offering type, equipment type and options, and may vary by country. Other restrictions may apply. Rates and offerings are subject to change, extension or withdrawal without notice.

IBM is not responsible for printing errors in this document that result in pricing or information inaccuracies.

All prices shown are IBM's United States suggested list prices and are subject to change without notice; reseller prices may vary.

IBM hardware products are manufactured from new parts, or new and serviceable used parts. Regardless, our warranty terms apply.

Any performance data contained in this document was determined in a controlled environment. Actual results may vary significantly and are dependent on many factors including system hardware configuration and software design and configuration. Some measurements quoted in this document may have been made on development-level systems. There is no guarantee these measurements will be the same on generally-available systems. Some measurements quoted in this document may have been estimated through extrapolation. Users of this document should verify the applicable data for their specific environment.