

Tiering in and Peering out Approaches to Data-Intensive Science

SOS XXIII, 28th March 2019

Mallikarjun (Arjun) Shankar, Ph.D. and Sudharshan Vazkhudai, Ph.D.

Oak Ridge Leadership Computing Facility National Center for Computational Science ORNL

ORNL is managed by UT-Battelle, LLC for the US Department of Energy



and Bevond





Science challenges for a smart supercomputer: What can a smart supercomputer do?

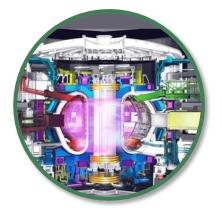


Leadership science challenges for a smart supercomputer:



Identifying Next-generation Materials

By training AI algorithms to predict material properties from experimental data, longstanding questions about material behavior at atomic scales could be answered for better batteries, more resilient building materials, and more efficient semiconductors.



Predicting Fusion Energy

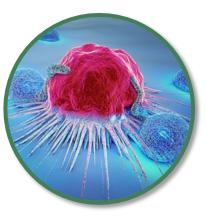
Predictive AI software is already helping scientists anticipate disruptions to the volatile plasmas inside experimental reactors. Summit's arrival allows researchers to take this work to the next level and further integrate AI with fusion technology.



National Laboratory

Deciphering High-energy Physics Data

With AI supercomputing, physicists can lean on machines to identify important pieces of information—data that's too massive for any single human to handle and that could change our understanding of the universe.

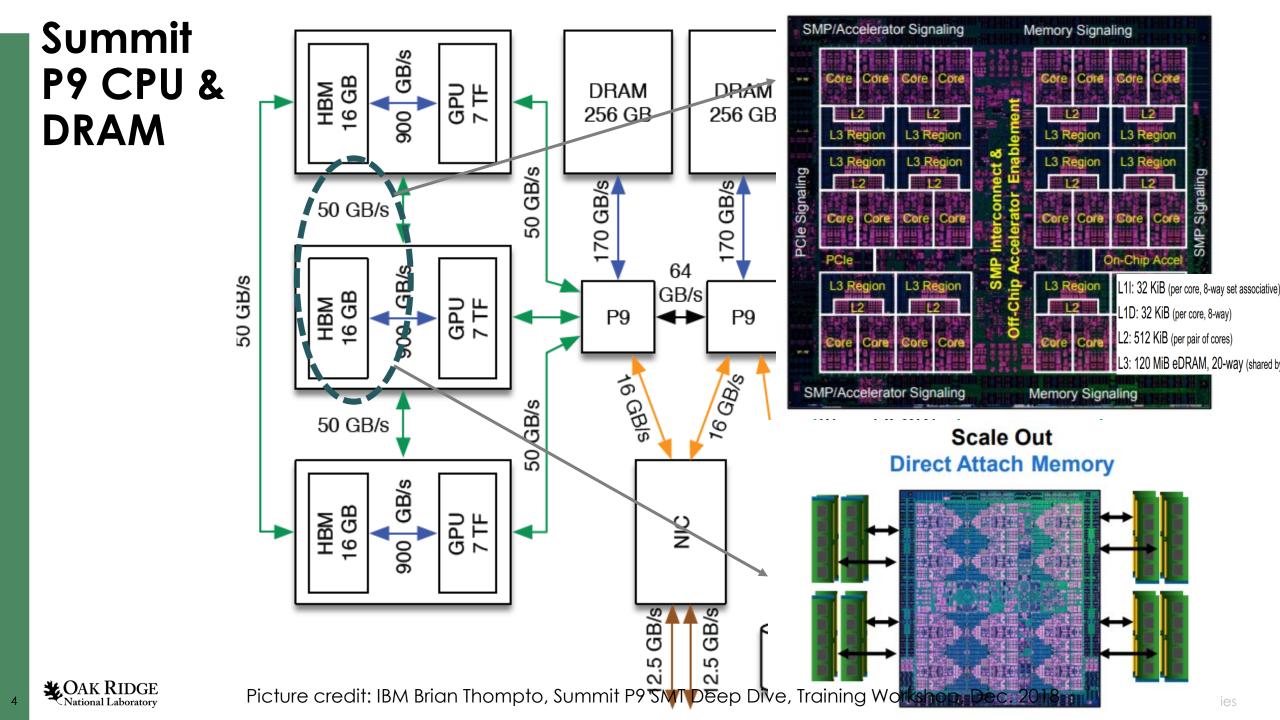


Combating Cancer

Through the development of scalable deep neural networks, scientists at the US Department of Energy and the National Cancer Institute are making strides in improving cancer diagnosis and treatment.

Our Tiered and Peered Data-Rich Landscape

- •We are in "Tiered" territory
 - -What is this doing for data science?
 - How does this inform our needs looking forward?
- "Peered" facilities for data-intensive science
 - Cross-facility data-intensive science
 - -The gap to address and our needs going forward?



GPU HBM

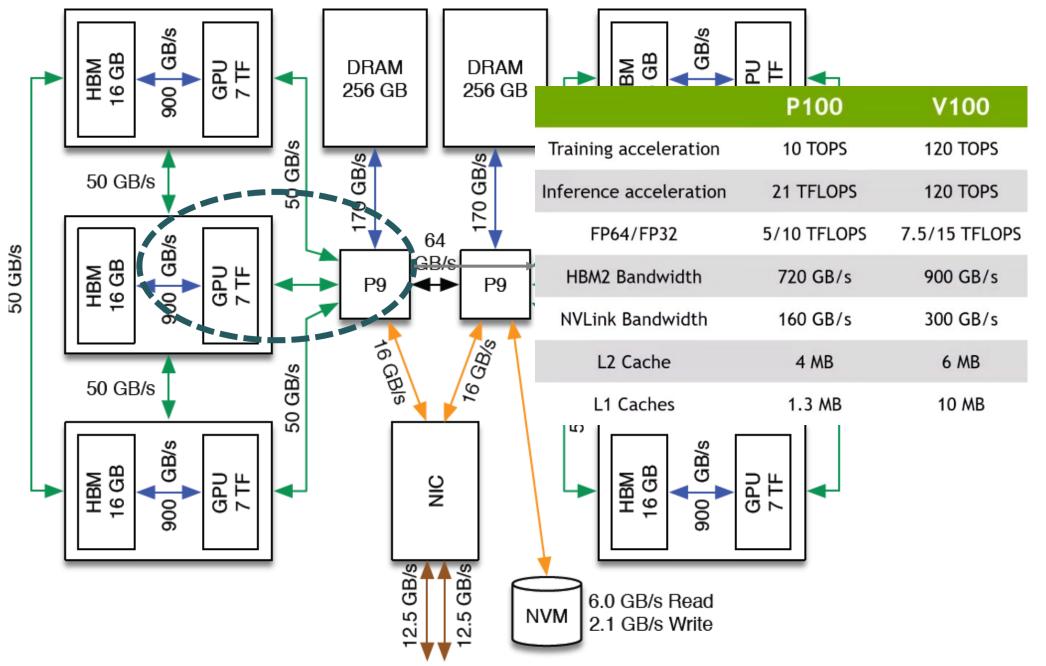
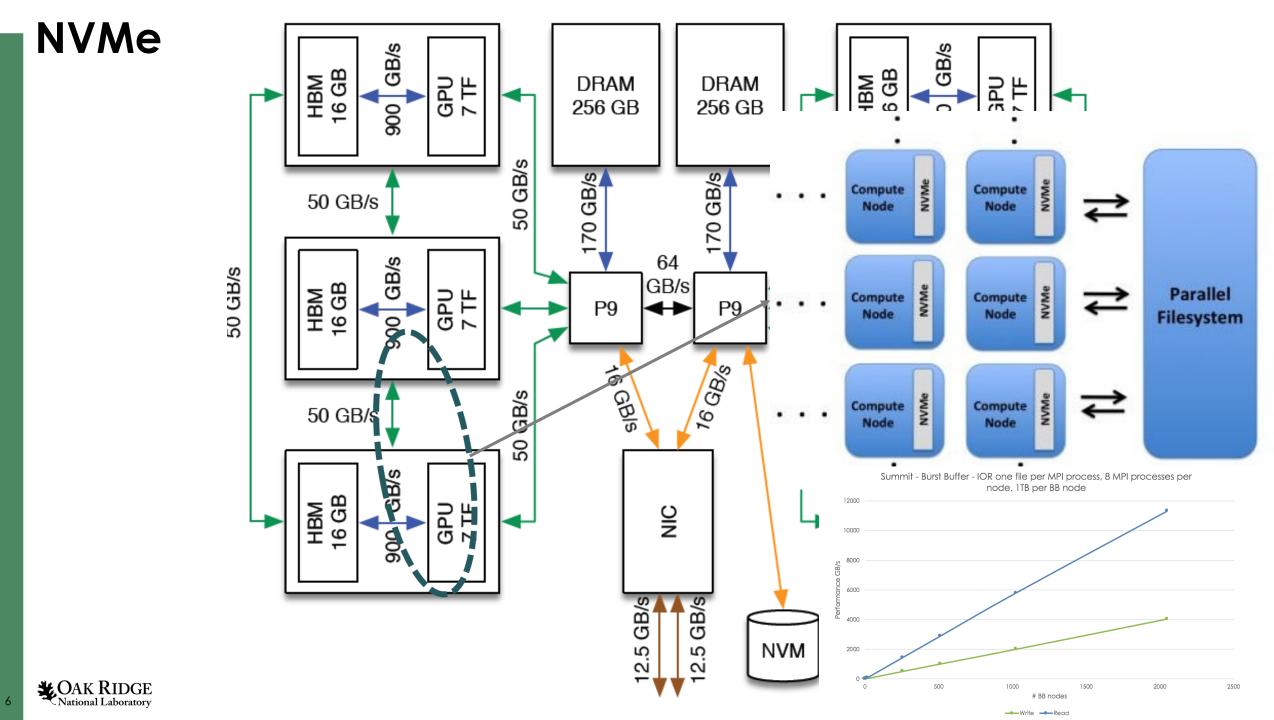
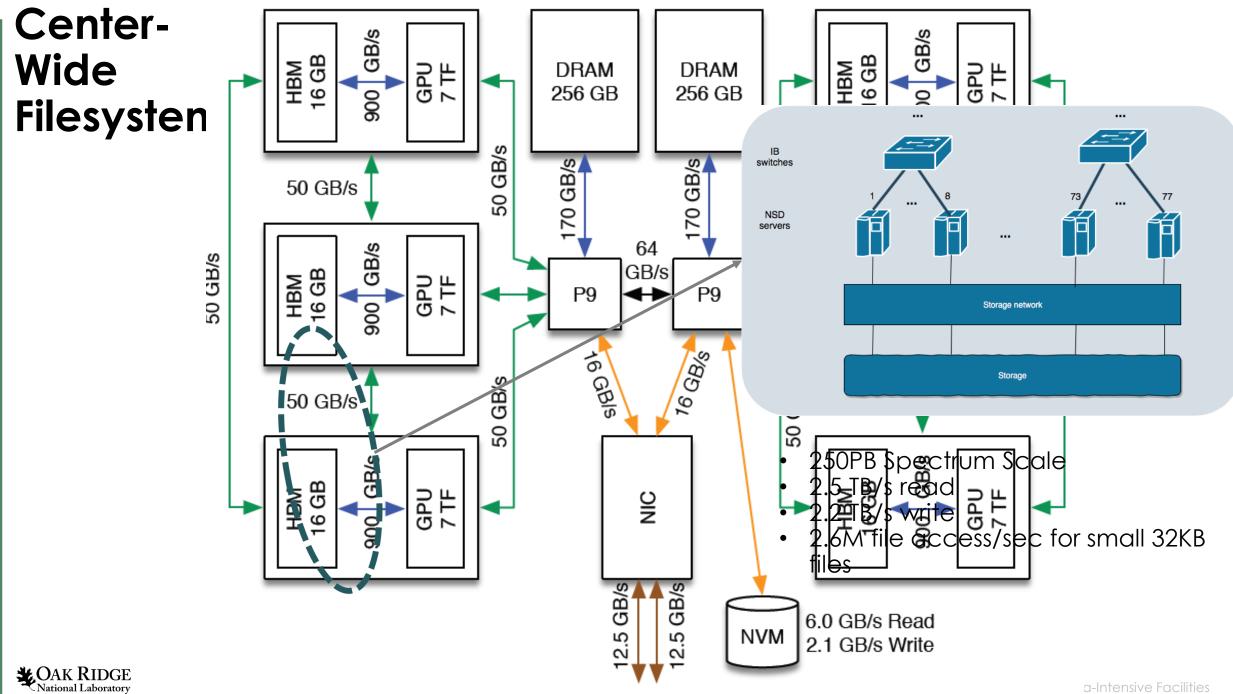


Table credit: NVIDIA, Jeff Larkin, Summit Iraining Workshop, Dec. 2018





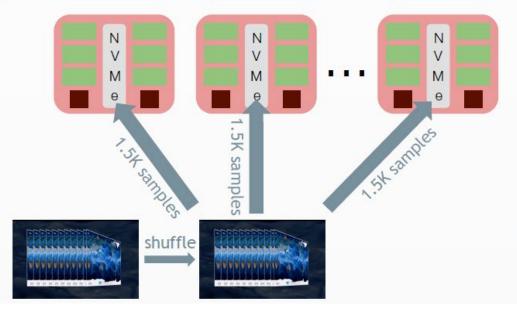


Example 1: Exascale Climate Application

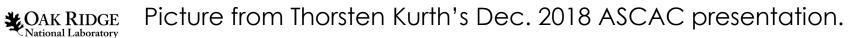
SC 18 Gordon Bell Award: Thorsten Kurth, Sean Treichler, Joshua Romero, Mayur Mudigonda, Nathan Luehr, Everett Phillips, Ankur Mahesh, Michael Matheson, Jack Deslippe, Massimiliano Fatica, Prabhat, Michael Houston; LBNL, NVIDIA, ORNL

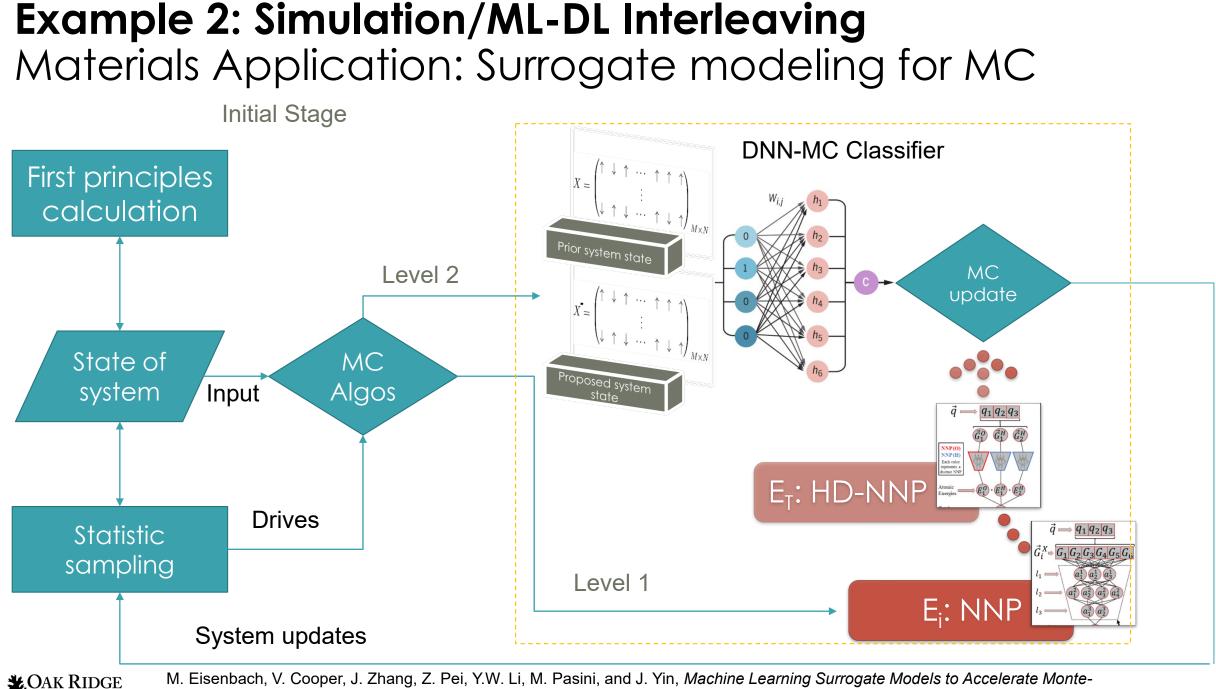
Data Staging

Dataset Size	Required BW (27K GPUs)	GPFS/LUSTRE	BurstBuffer	NVM/e or DRAM
20 TB (~63K samples)	3.8 TB/s	~400 GB/s	~2 TB/s	~26 TB/s



- 250 training samples/GPU (~15 GB), sample w/ replacement
- each file will be read at most once from FS
- files shared between nodes via MPI (mpi4py)





Carlo Calculation, Bulletin of the American Physical Society, March 2019

National Laboratory

Tiering and Peering for Data-Intensive Facilities

Data Science and Learning on the Memory Hierarchy

CORAL 2 Benchmark Suite	Description	
Big Data Analytic Suite (BDAS)	PCA, K-Means, and SVM (based on pbdR)	
Deep Learning Suite (DLS)	CANDLE, CNN, RNN, and ResNet-50 (distributed memory)	
Deep Learning Codes (CNN; ResNet50;) excel here with NVM and GPUs enabling tensor operations.		50 GB/s 50

CAK RIDGE

Traditional Node: PCA, K-Means, etc. excel due to the node's memory, CPU, and on-chip bandwidth

Programming with Big Data in R

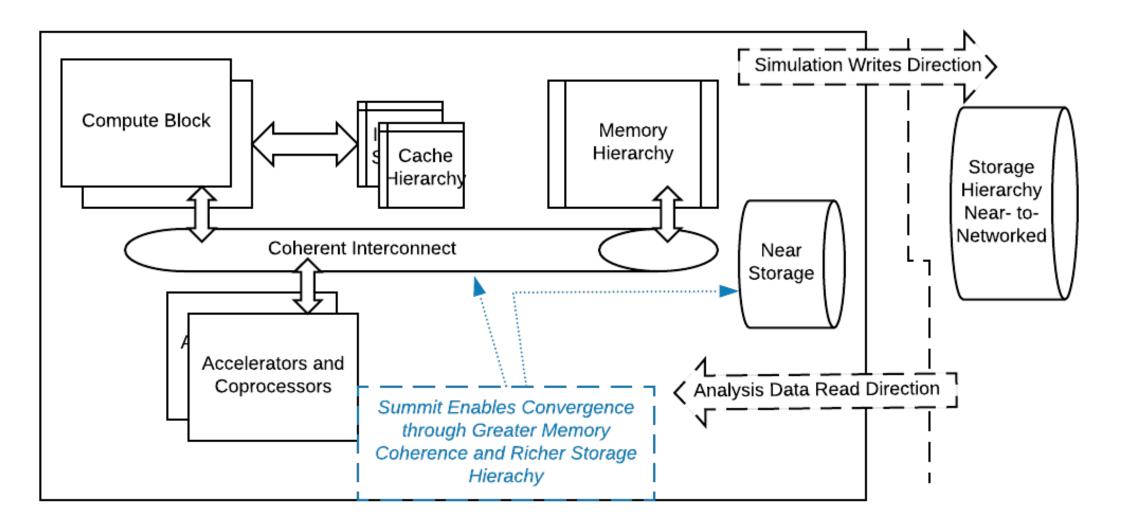
Code suites are in the CORAL (Collaboration of Oak Ridge, Argonne, Livermore laboratories data benchmark suite: https://asc.llnl.gov/coral-2-benchmark**\$0**

Accommodating the Data-Science Landscape

11

Model-Driven and Simulation-Heavy Workloads	Data-Driven and Analytics (ML/DL) Workloads
Single Center-Wide Namespace	Read-Intensive Partitioned Namespace OK
Write Fraction (~15%) of HBM Rapidly (~5 Minutes)	High IOPs
Large Capacity Retention	Large Capacity Retention
N-N, N-1 write access patterns	Small, Random Reads

Tiering has moved us towards: "Convergence"





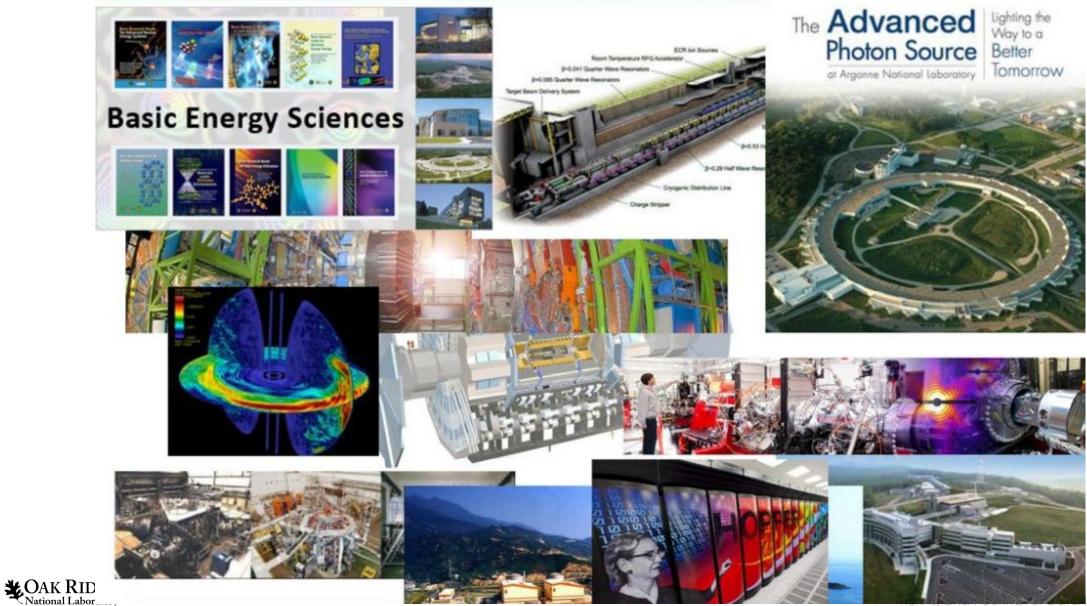
A rich tiered data hierarchy is a necessity given our **technology** (bandwidth, latency) and **cost** (per bit) landscape!

Tools to help the user to seamlessly process data across node-local storage, and node-local and center-wide storage are needed.

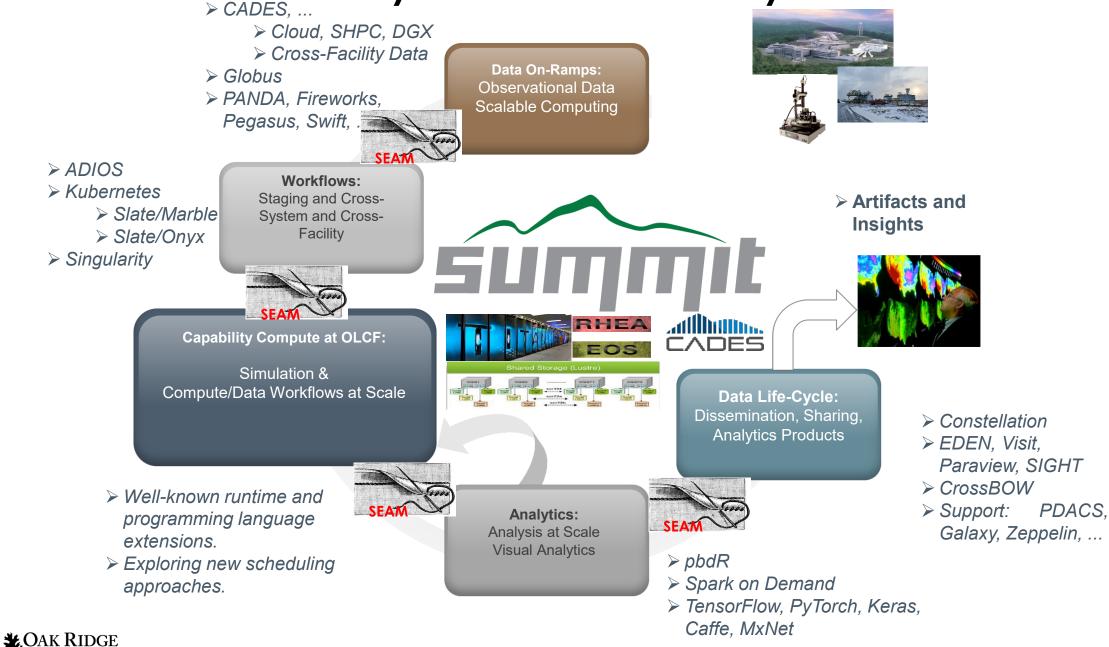
Shameless plug: ORNL is working on a couple - Spectral and Symphony FS.



HPC Cloud-Edge-Fogging up the data-story...



Scalable Cross-Facility Workflow and Ecosystem for Data Science



15

National Laboratory

Tiering and Peering for Data-Intensive Facilities

To Conclude..

1. Tiering: Need tools to support data processing across heterogeneous compute and memory hierarchy, node-local storage, and node-proximal and center-wide storage.

2. Peering: Need tools to support distributed data read/process/write abstractions across data-stores in data-intensive facilities.

#1 and #2 call for similar conceptual abstractions of programming primitives. The implementation will differ based on data affinity, currency, and latency awareness. Both could use an underlying data handling layer/API.



Thank You!

Acknowledgements: OLCF/NCCS Colleagues (Junqi Yin, Jack Wells, Chris Zimmer); IBM; NVIDIA; Mellanox

This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory (ORNL). ORNL is operated by UT-Battelle, LLC, for the U.S. Department of Energy under contract DE-AC05-00OR22725.

The United States Government retains a non-exclusive, paid-up, irrevocable, worldwide license to publish or reproduce the published form of this manuscript, or allow others to do so, for United States Government purposes. In addition, this work is supported by the Laboratory Director R&D Fund.



Peering Reinterpreted





Observational Facility This Hildub by Unknown Author Vicensed under <u>CC BY-SA</u>

is licensed under <u>CC BY-SA</u>

Computational Facility