

# Data Science at OLCF

Bronson Messer  
Scientific Computing Group  
Oak Ridge Leadership Computing Facility  
Oak Ridge National Laboratory

Mallikarjam “Arjun” Shankar  
Group Leader – Advanced Data and Workflows Group  
Oak Ridge Leadership Computing Facility  
Oak Ridge National Laboratory

# OLCF Data/Learning Strategy & Tactics

Applications

1. Engage with applications
  - Summit Early Science Applications (e.g., CANDLE)
  - INCITE projects (e.g., Co-evolutionary Networks: From Genome to 3D Proteome, Jacobson, et al.)
  - Directors Discretionary projects (e.g., Fusion RNN, MiNerva)

Algorithms

2. Create leadership-class analytics capabilities
  - Leadership analytics (e.g., Frameworks: pbdR, TensorFlow + Horovod)
  - Algorithms requiring scale (e.g., non-negative matrix factorization)

Infrastructure

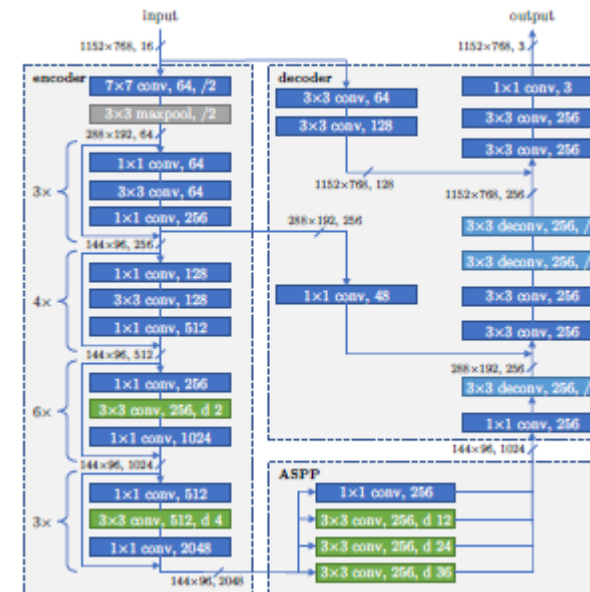
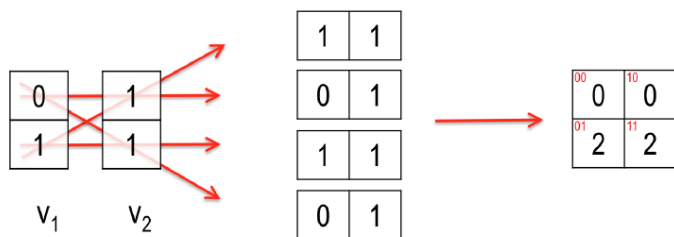
3. Enable infrastructure for analytics/AI and data-intensive facilities
  - Workflows to include data from observations for analysis within OLCF
  - Analytics enabling technologies (e.g., container deployments for rapidly changing DL/ML frameworks, analytics notebooks, etc.)

# Applications Supported through DD/ALCC: Selected Machine Learning Projects on Titan: 2016-2017

Program	PI	PI/Employer	Project Name	Allocation (Titan Core-hrs)
ALCC	Robert Patton	ORNL	Discovering Optimal Deep Learning and Neuromorphic Network Structures Using Evolutionary Approaches on High Performance Computers	75,000,000
ALCC	Gabriel Perdue	FNAL	Large Scale Deep Neural Network Optimization for Neutrino Physics	58,000,000
ALCC	Gregory Laskowski	GE	High-Fidelity Simulations of Gas Turbine Stages for Model Development Using Machine Learning	30,000,000
ALCC	Efthimios Kaxiras	Harvard U.	High-Throughput Screening and Machine Learning for Predicting Catalyst Structure and Designing Effective Catalysts	17,500,000
ALCC	Georgia Tourassi	ORNL	CANDLE Treatment Strategy Challenge for Deep Learning Enabled Cancer Surveillance	10,000,000
DD	Abhinav Vishnu	PNNL	Machine Learning on Extreme Scale GPU Systems	3,500,000
DD	J. Travis Johnston	ORNL	Surrogate Based Modeling for Deep Learning Hyper-parameter Optimization	3,500,000
DD	Robert Patton	ORNL	Scalable Deep Learning Systems for Exascale Data Analysis	6,500,000
DD	William M. Tang	PPPL	Big Data Machine Learning for Fusion Energy Applications	3,000,000
DD	Catherine Schuman	ORNL	Scalable Neuromorphic Simulators: High and Low Level	5,000,000
DD	Boram Yoon	LANL	Artificial Intelligence for Collider Physics	2,000,000
DD	Jean-Roch Vimant	Caltech	HEP Deep Learning	2,000,000
DD	Arvind Ramanathan	ORNL	ECPC Cancer Distributed Learning Environment	1,500,000
DD	John Cavazos	U. Delaware	Large-Scale Distributed and Deep Learning of Structured Graph Data for Real-Time Program Analysis	1,000,000
DD	Abhinav Vishnu	PNNL	Machine Learning on Extreme Scale GPU Systems	1,000,000
DD	Gabriel Perdue	FNAL	MACHINE Learning for MINERVA	1,000,000
		<b>TOTAL</b>		<b>220,500,000</b>

 - Highlighted rows are Algorithm and Infrastructure Examples;  
Rest are Primarily Science Applications

# Gordon Bell Prizes in 2018: Peak Performance on Summit



## Attacking the Opioid Epidemic: Determining the Epistatic and Pleiotropic Genetic Architectures for Chronic Pain and Opioid Addiction

Wayne Joubert, Deborah Weighill, David Kainer, Sharlee Climer, Amy Justice, Kjersten Fagnan, Daniel Jacobson

## Exascale Deep Learning for Climate Analytics

Thorsten Kurth, Sean Treichler, Joshua Romero, Mayur Mudigonda, Nathan Luehr, Everett Phillips, Ankur Mahesh, Michael Matheson, Jack Deslippe, Massimiliano Fatica, Prabhat, Michael Houston,

# Methods: Leadership Data Analytics Capabilities

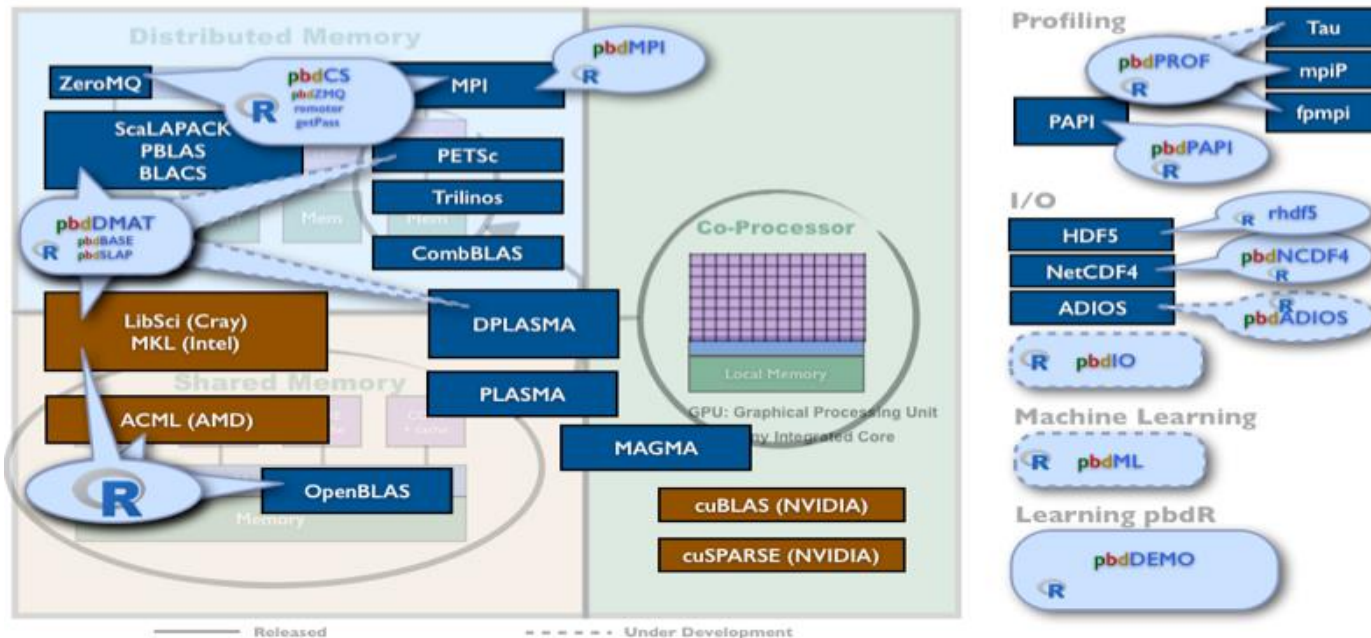


- Engage parallel math libraries at scale
- R language unchanged
- New distributed concepts
- New profiling capabilities
- New interactive SPMD parallel
- In-situ distributed capability
- In-situ staging capability via ADIOS

<http://pbdr.org>

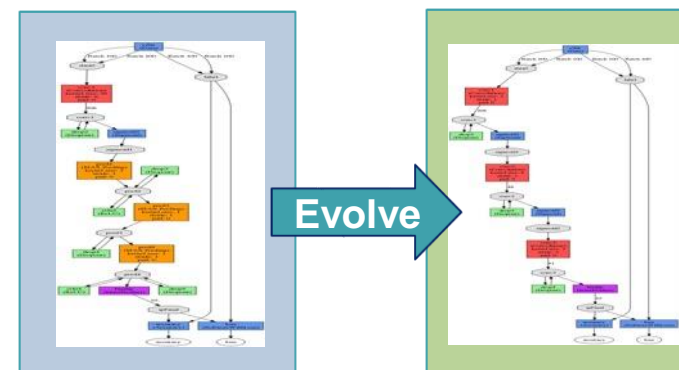
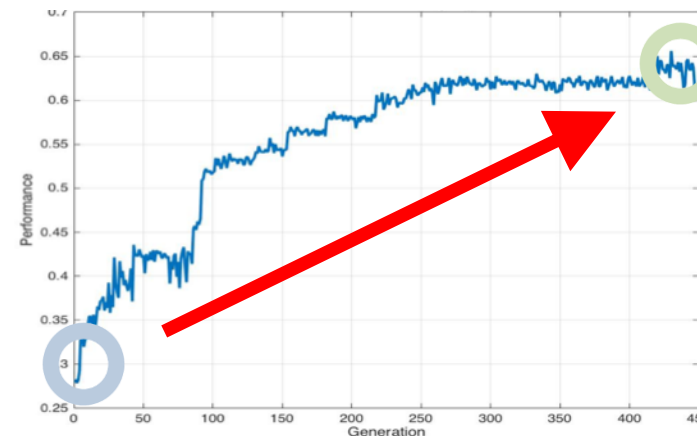
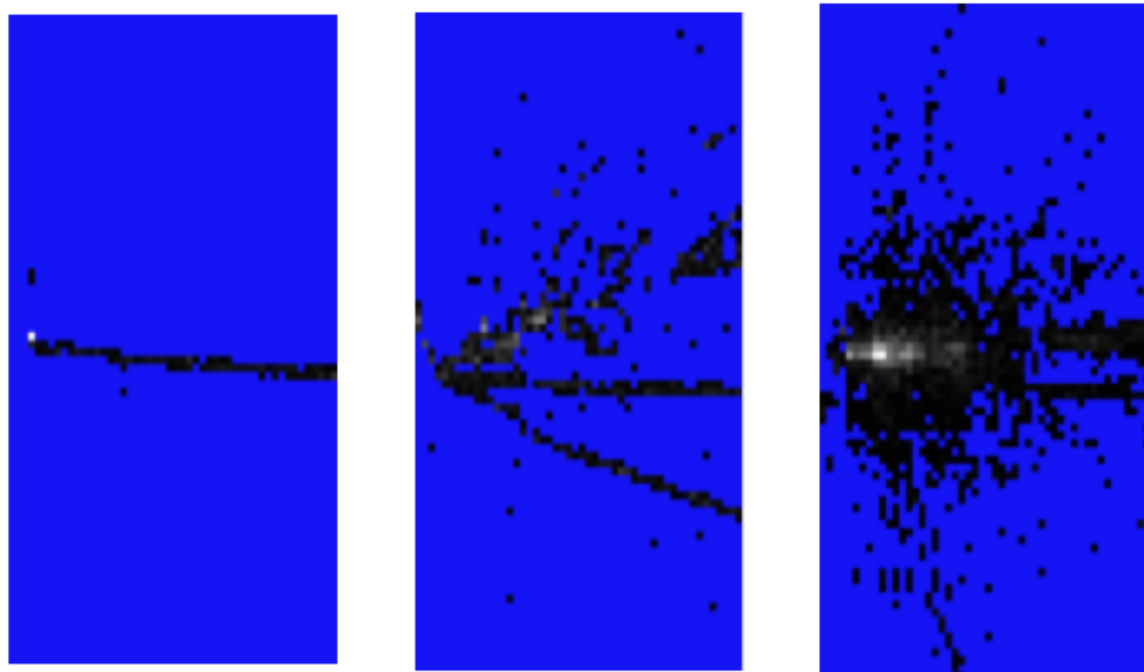
- Significantly improved throughput compared to, e.g., Apache Spark
  - “PCA of a 134 GB matrix: ‘hours on . . . Apache Spark, . . . less than a minute using R.’”
  - June 2016, HPCWire

HPC libraries and their R/pbdR connections



# Scaling Deep Learning for Science with ORNL's MENNDL

ORNL-designed algorithm leverages Titan to create high-performing neural networks



An image generated from neutrino scattering data captured by the MINERvA detector at Fermi National Accelerator Laboratory. Researchers are using MENNDL and the Titan supercomputer to generate deep neural networks that can classify high-energy physics data and improve the efficiencies of measurements.

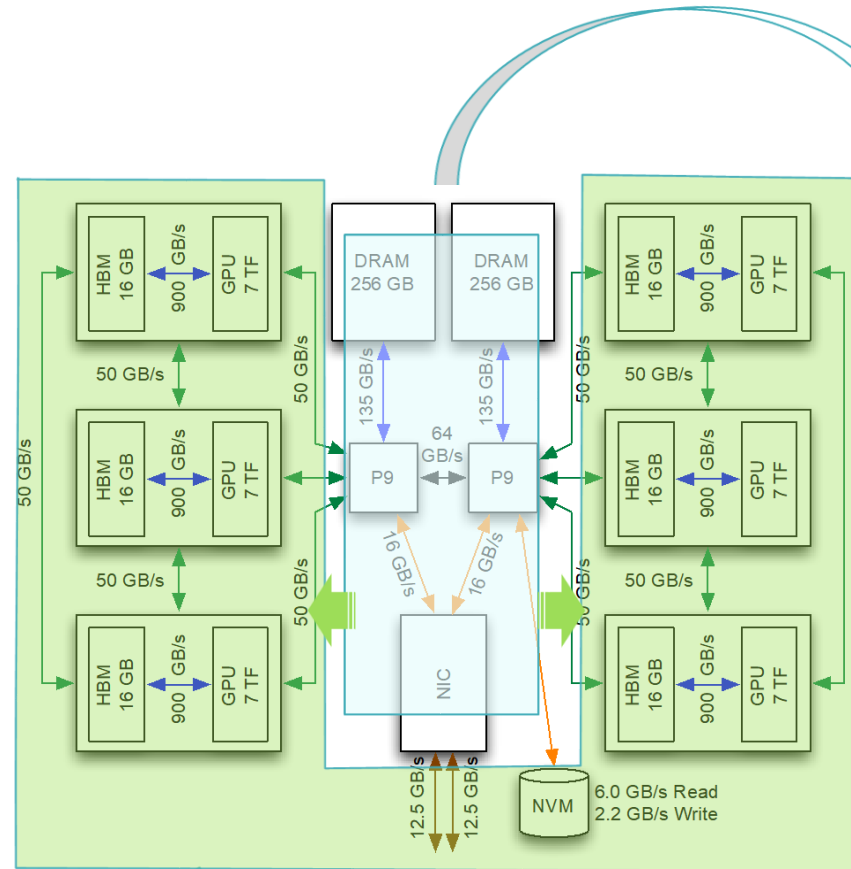
ORNL Data Analytics Group used Titan to develop an evolutionary algorithm to search for optimal hyper-parameters and topologies for ML networks.

Steven R. Young, Derek C. Rose, Travis Johnston, William T. Heller, Thomas P. Karnowski, Thomas E. Potok, Robert M. Patton, Gabriel Perdue, and Jonathan Miller, "Evolving Deep Networks Using HPC." In *Proceedings of the Machine Learning on HPC Environments*. Paper presented at *The International Conference for High Performance Computing, Networking, Storage and Analysis*, Denver, Colorado (November 2017), doi: [10.1145/3146347.3146355](https://doi.org/10.1145/3146347.3146355).

# Key Data Science and Learning Methods

CORAL 2 Benchmark Suite	Components
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.

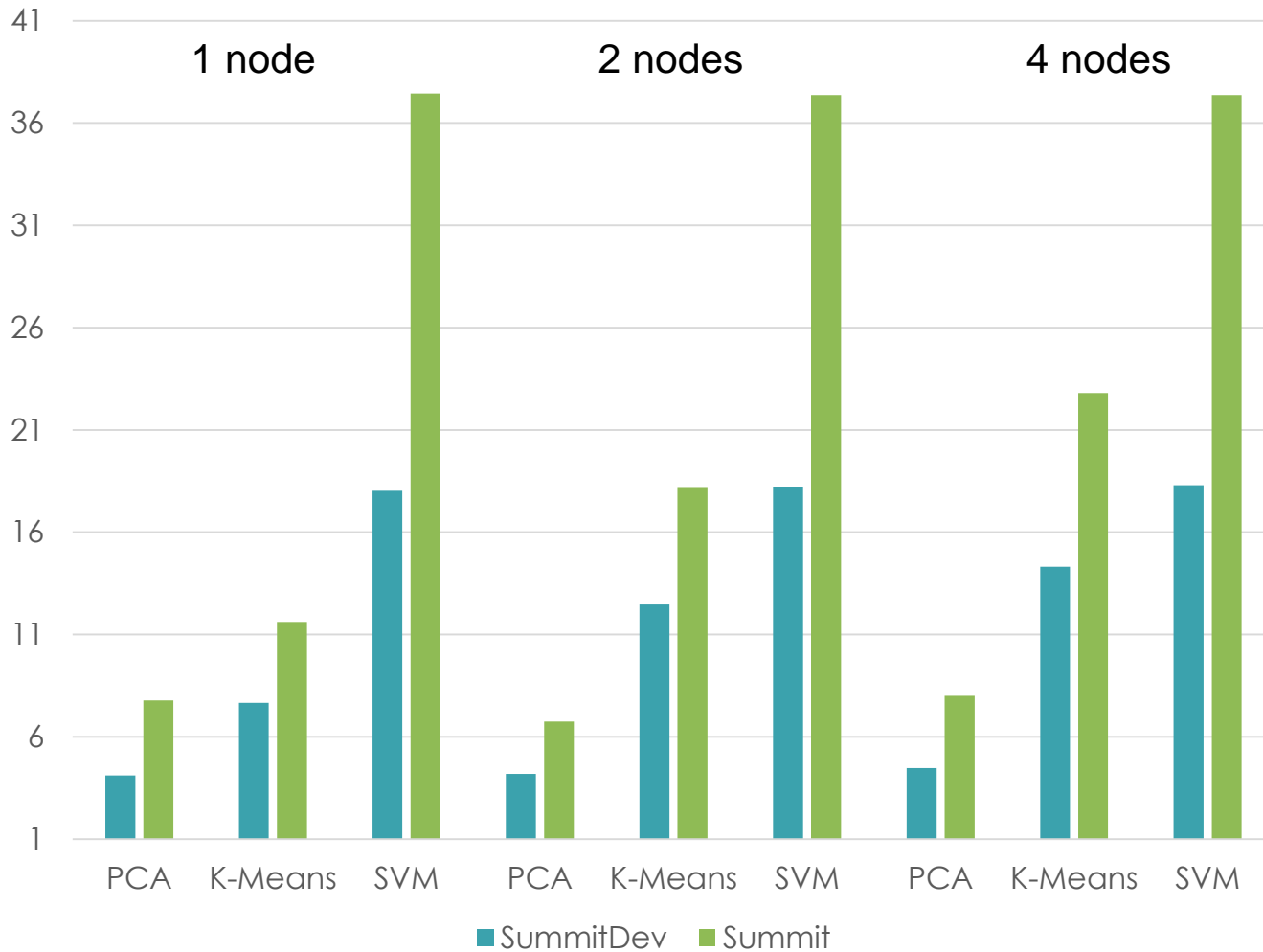


PCA, K-Means, etc. excel on “traditional HW” part of the node due to the node’s memory, CPU, and on-chip bandwidth

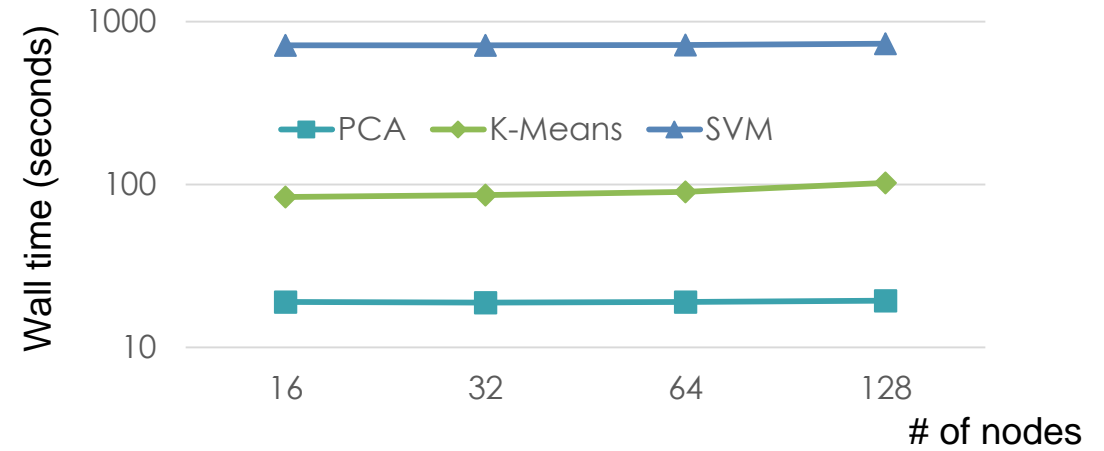
Code suites are in the CORAL2 (Collaboration of Oak Ridge, Argonne, Livermore laboratories) benchmark suite: <https://asc.llnl.gov/coral-2-benchmarks/>

# Big Data Analytic Suite

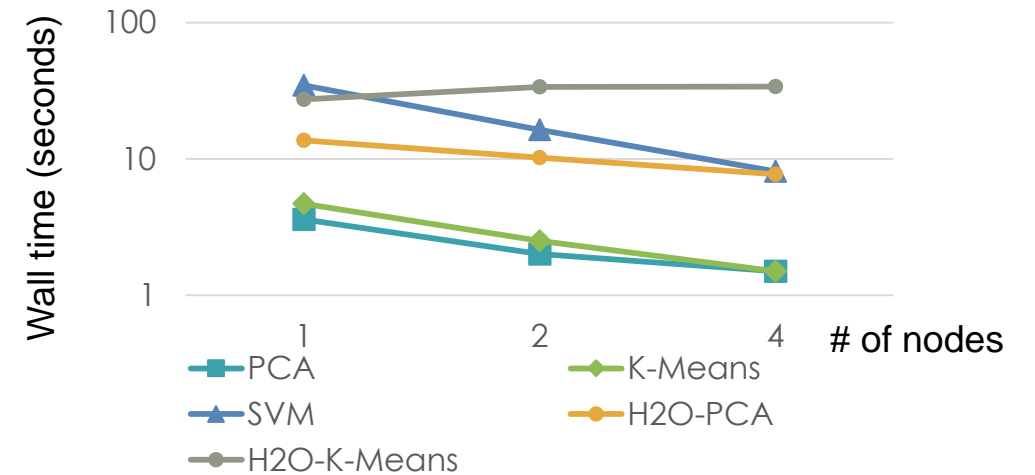
Speedup Over Titan Baseline for CORAL-2  
Big Data Benchmarks (based on pbdR)



Weak Scaling of Data Benchmarks on Titan



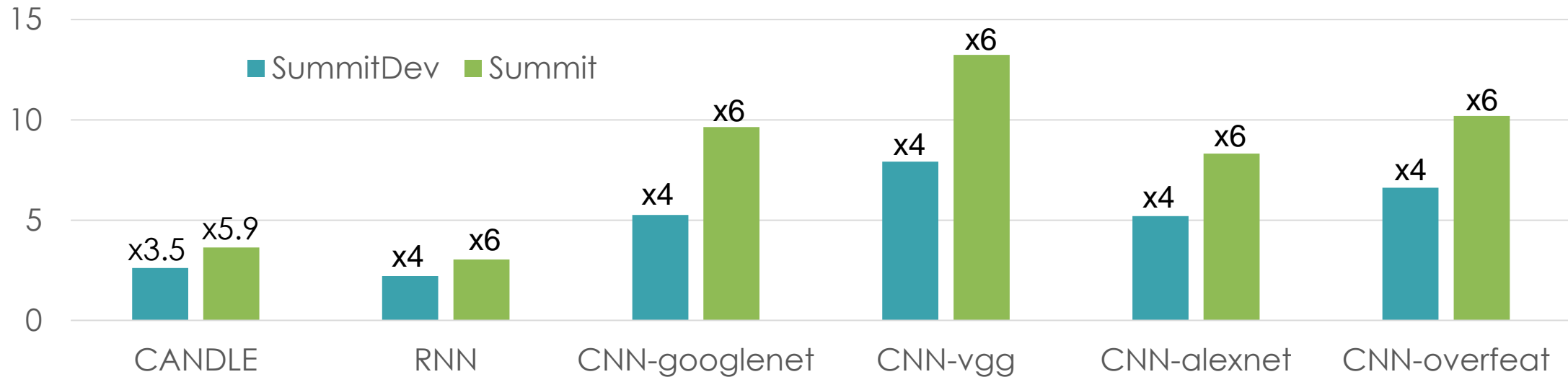
Strong Scaling of Data Benchmarks on SummitDev



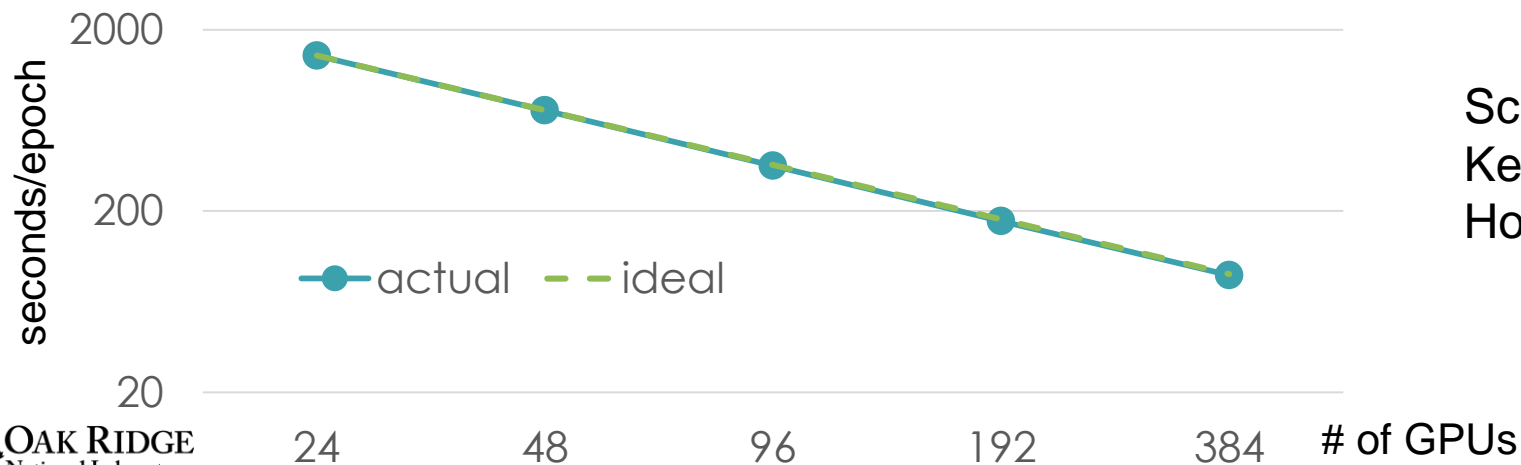


# Deep Learning Suite

Speedup Over Titan Baseline for CORAL-2 Deep Learning Benchmarks

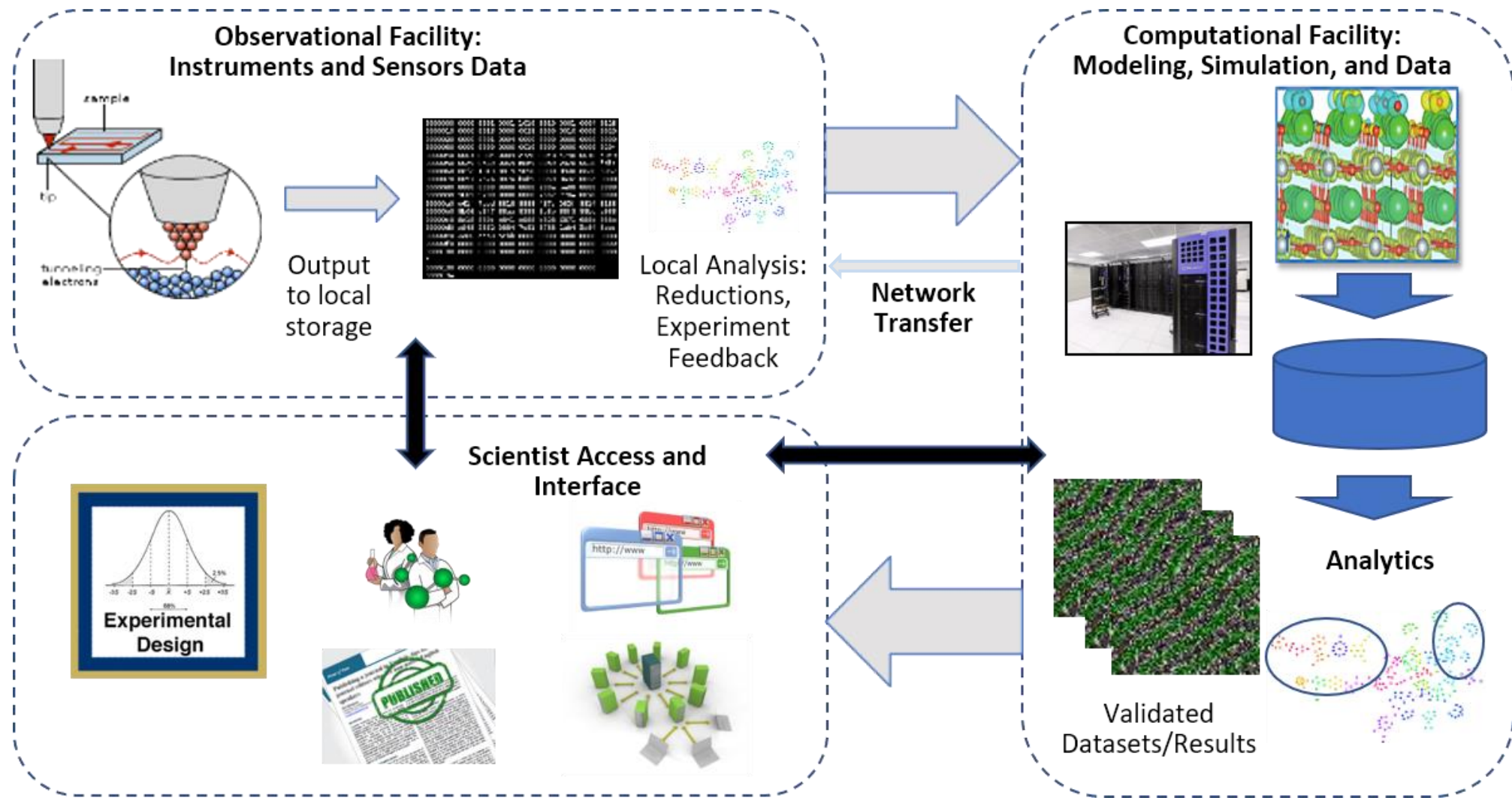


Strong Scaling of ResNet-50 on Summit



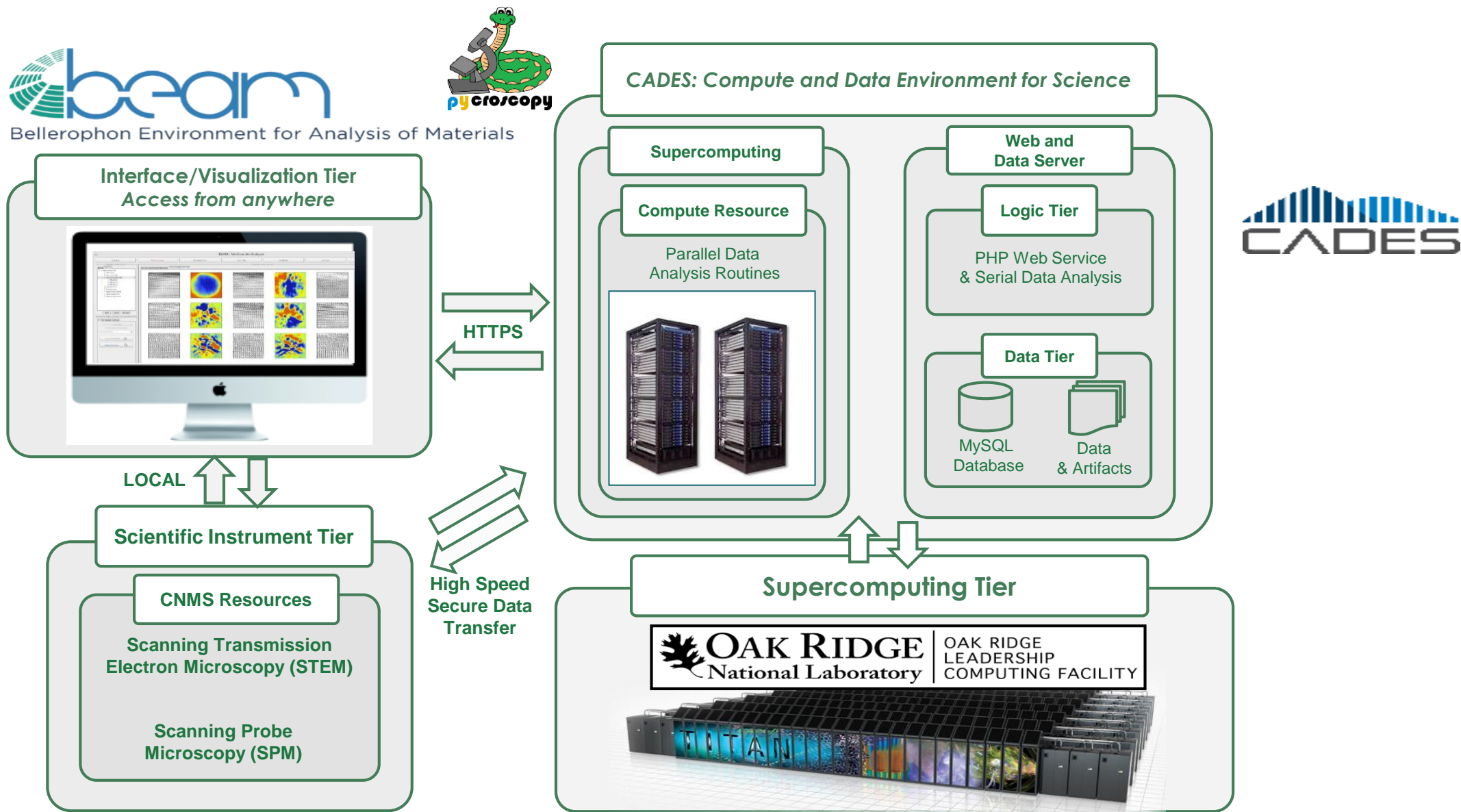
Scaling of Resnet-50 based on Keras (Tensorflow backend) and Horovod on ImageNet data

# Infrastructure - Cross-Facility Design Pattern



From: **Policy Considerations when Federating Facilities for Experimental and Observational Data Analysis**, Mallikarjun (Arjun) Shankar, Suhas Somnath, Sadaf Alam, Derek Feichtinger, Leonardo Sala, and Jack Wells, (2018, Submitted Book Chapter)

# CADES runs BEAM and Pycroscopy for SNS and CNMS



# Towards Data Service Offerings and Easier Data Access Across Facilities

- Categories of Data Services
  - **Type 1** data repository program for “**data-only**” projects.
  - **Type 2** data services program for user communities.
  - **Type 3** computational and data science end station program.
- “DataFed” prototype to enable federated data access across facilities (currently being tested)

