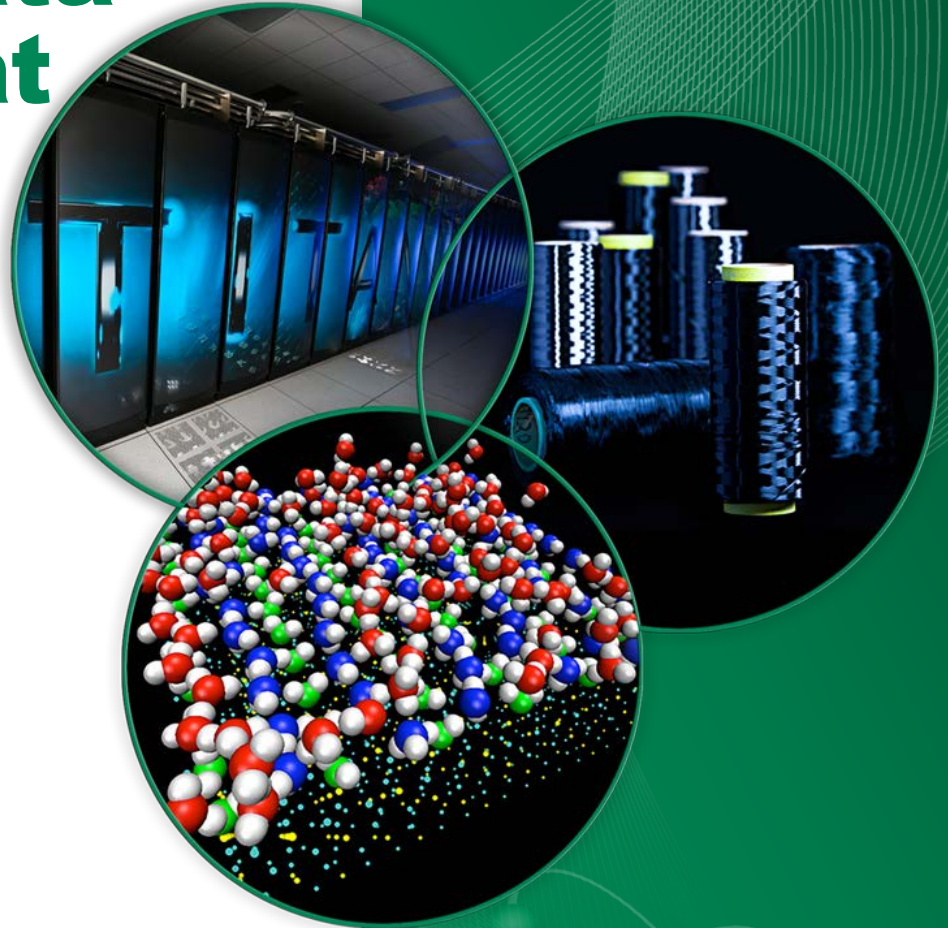# The Future of Data and Computing at ORNL

Barney Maccabe

April 29, 2013

# DOE is a world leader in HPC
## *- leaders in architecting, acquiring and deploying computers used as instruments of discovery*

- A record of leadership in computing performance

- Thousands of users from government labs, universities, and industry

- Forefront computing facilities
  - Used to solve mission problems in nuclear defense, science, and engineering
  - Enabling prize-winning science (e.g., Nobel), engineering solutions, and thousands of publications

**TOP500**
November 2012



Titan at ORNL
(#1, 17+ PF)

Sequoia at LLNL
(#2, 16+ PF)

Mira at ANL
(#4, 8+ PF)

Cielo at LANL/SNL
( #18, 1+PF)

Hopper at LBNL
(#19, 1+PF)

**OAK RIDGE NATIONAL LABORATORY**
MANAGED BY UT-BATTELLE FOR THE U.S. DEPARTMENT OF ENERGY

# Scientific Computing as an Instrument
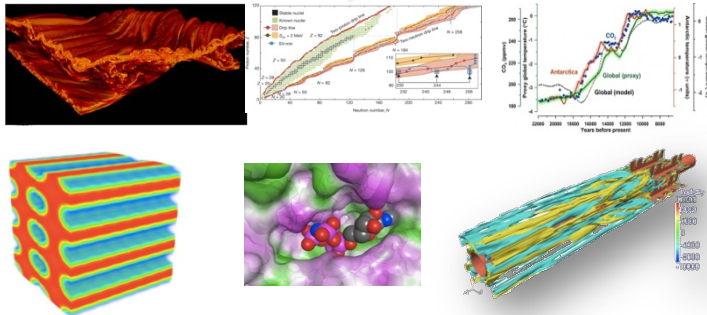
$$\left(-\nabla^2 + V\right)\Psi = E\Psi$$

$$\Psi = -\left(-\nabla^2 - E\right)^{-1}V\Psi$$

$$= -G*(V\Psi)$$

$$(G*f)(r) = \int ds \frac{e^{k|r\ s|}}{4\pi|r-s|}f(s) \text{ in 3D }; k^2 = -E$$
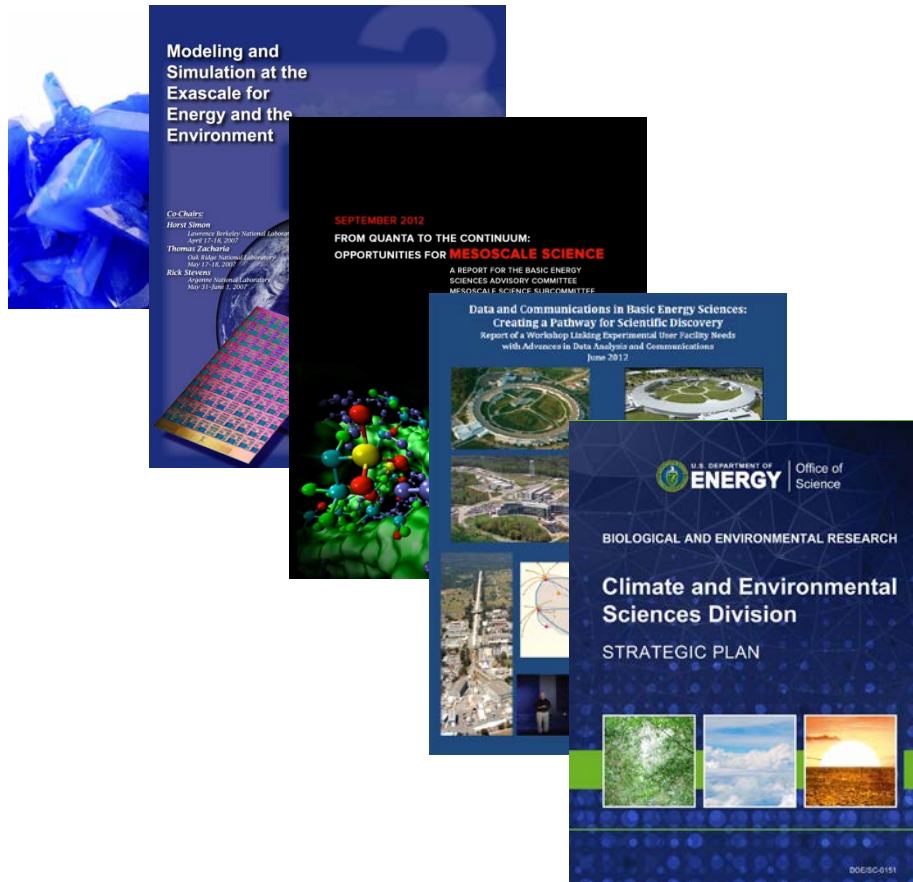
- Over the past few decades we have developed world-class facilities for scientific simulation

- These "simulation instruments" are now integral to theory based scientific discovery

- Titan will support a broad range of theoretic studies in Astrophysics, Climate Science, Materials Science, Nuclear Physics, Fusion, etc.

OAK RIDGE NATIONAL LABORATORY
MANAGED BY UT-BATTELLE FOR THE U.S. DEPARTMENT OF ENERGY

# Data Infrastructure and Scientific Discovery
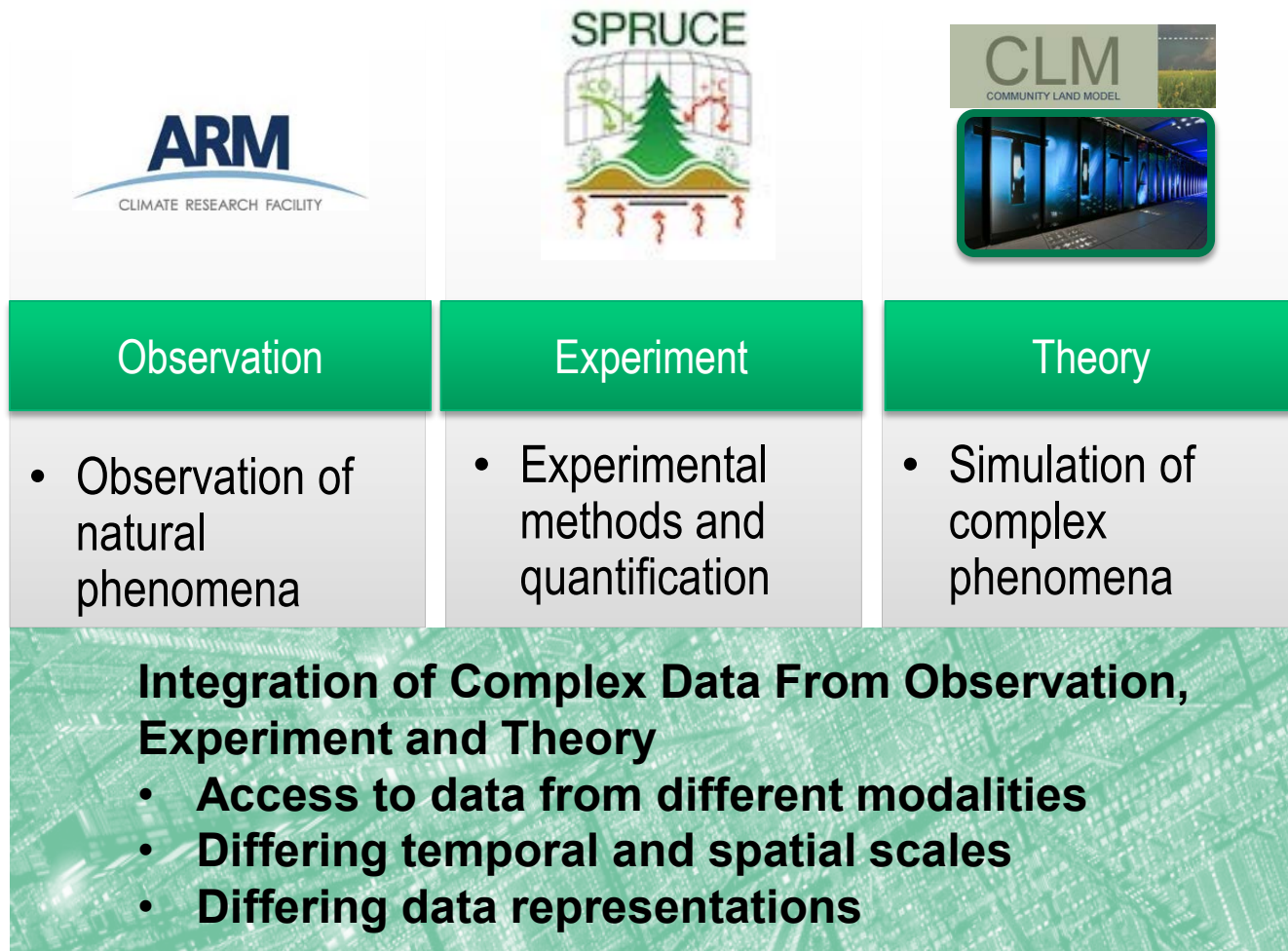


- To date, our primary focus in has been simulation

- However, the rate of scientific progress is increasingly dependent on the ability to efficiently capture, integrate, analyze, and steward large volumes of diverse data

- Increasing data volume, variety, and velocity require a new environment for scientific discovery

- Major facilities and research programs across the Office of Science must develop adequate infrastructure to support this new environment

# Data Intensive Computing as Integrative Infrastructure

| Observation | Experiment | Theory |
|---|---|---|
| • Observation of natural phenomena | • Experimental methods and quantification | • Simulation of complex phenomena |

**Integration of Complex Data From Observation, Experiment and Theory**
- **Access to data from different modalities**
- **Differing temporal and spatial scales**
- **Differing data representations**

# Office of Science Response To Data Intensive Challenges

- The ASCR (Advanced Scientific Computing Research) program is gearing up to address these needs

  *SC 2013 Budget request*

  – Computer science research

    "FY 2013 supports new research efforts to address the challenges of data-intensive science with focus on full data lifecycle management and analysis for the massive data from DOE scientific user facilities."

  – Computational Partnerships

    "New research efforts will engage partners across the Office of Science to address the data-intensive science challenges at the science application level."

  – Next Generation Networking

    "FY 2013 supports new research efforts to address the data-intensive science challenges facing Request scientific communities using unique DOE facilities and engaging in large-scale collaborations."

- **What is missing?**

  – **Facilities to support the data challenges facing the Office of Science**

OAK RIDGE NATIONAL LABORATORY
MANAGED BY UT-BATTELLE FOR THE U.S. DEPARTMENT OF ENERGY

# A "Virtual Data Facility" for Office of Science

- The Virtual Data Facility will provide these capabilities
  - Rich data analysis environment - tightly coupling compute and data storage
  - Coupling simulation, experiment, and observation data
    - Model improvement, validation, steering, site selection, etc.
  - A flexible compute and data environment that can be tailored to specific needs
  - Cataloging and long-term stewardship of scientific datasets
  - Long-term (many year) allocations for major facilities and projects
- The Virtual Data Facility will be built as an extension of the three compute facilities at Argonne, Berkeley, and Oak Ridge
  - Core ASCR facilities (OLCF, ALCF, and NERSC) have already been designed to meet the needs of computational science (simulation) workloads
  - These facilities will be extended to meet the data science needs
    - Leverages facilities investments and staff expertise at the OLCF, ALCF, NERSC, and ESnet
- The Virtual Data Facility will  create a new environment for scientific discovery for the office of science (ASCR, BER, BES, NP, FES, HEP, NE, etc.)

# The Virtual Data Facility

- **Integrated Data Science Infrastructure across ASCR Facilities in support of major projects and facilities**



- Data Science Liaisons
- Flexible Software Compositions
- Large-scale data resources
  - 100 Petabytes growing to 5 Exabytes
  - 3 Petaops/s growing to 300 Petaops/s
- Long term access & stewardship

# Data Science and Compute Strategy at ORNL

- Consolidate resources and develop a flexible and elastic infrastructure to satisfy needs of existing and future computing and data projects

- Establish a Compute & Data Environment for Science (CADES)

- Provide competitive advantage in R&D and for national data centers (DSF)

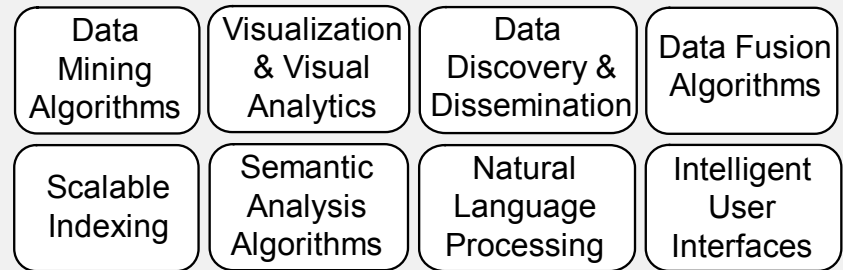Compute and Data Environment for Science targets full range of needs

| Single PI Projects | Centers | User Facilities | Multi-Institutional Projects |
|---|---|---|---|

## ORNL Compute and Data Environment for Science

### *Compute and Data Science Liaisons*

**Compute and Data Platforms**     **Tools, Software**

| Data Analysis Systems | File Systems | Structured Storage Databases | Mid-Range Compute | | Data Mining Algorithms | Visualization & Visual Analytics | Data Discovery & Dissemination | Data Fusion Algorithms |
|---|---|---|---|---|---|---|---|---|
| Semi-Structured Storage | Ubiquitous Storage (Cloud) | Archival Storage | Cloud Compute | | Scalable Indexing | Semantic Analysis Algorithms | Natural Language Processing | Intelligent User Interfaces |

Supply Common Computing and Data Needs

# CADES Infrastructure

## People

**Matrix staff with expertise in all areas of compute and data science**

- Domain-specific compute and data science

- Algorithms, data analysis, and visualization

- Compute and data system architecture

- Operations

## Flexible and Elastic Resources

- Flexible to meet requirements of a broad set of initiatives
  - Performance, scalability, manageability, security
  - Compute: From enterprise servers and clusters to specialized systems
  - Storage: From network attached storage and parallel file systems to archive
  - Software: From commercial or community to custom packages

- Elastic to meet on-demand compute and storage requirements for data-intensive workloads
  - Projects no longer constrained by a fixed system resource

# Providing Expertise in Compute & Data For Mission

**CADES is like a hub -- it shares data infrastructure and Compute & Data Science capabilities with and among many projects**

**Projects with core competencies in Compute and Data Science can share these expertise via the HUB**

**CADES is delivering to initiatives today and provides the necessary capabilities required by the SC Data Science Facility in the future**
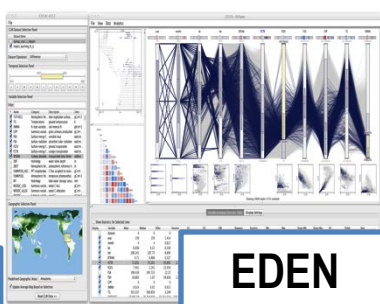
# CADES – Delivering to Initiatives Today

## Scalable tools and data infrastructure for climate science

- Data discovery and dissemination tools
- Scalable data analytics and data fusion
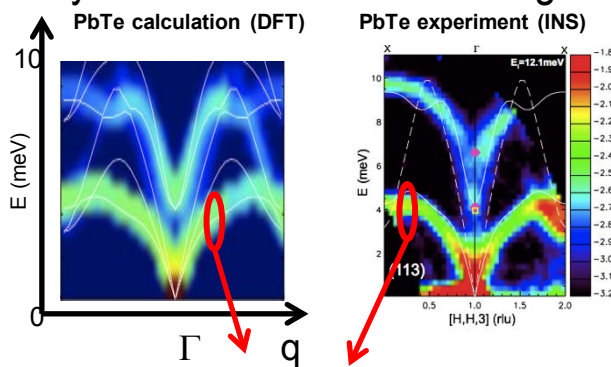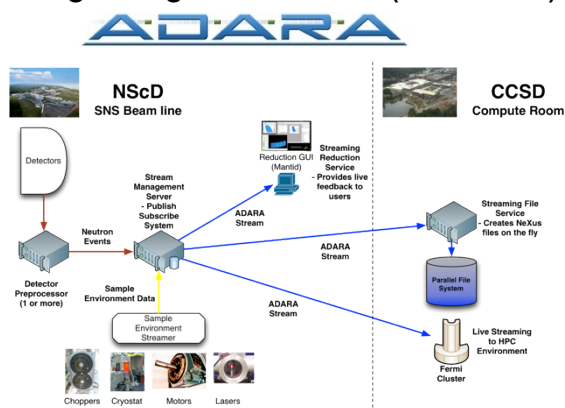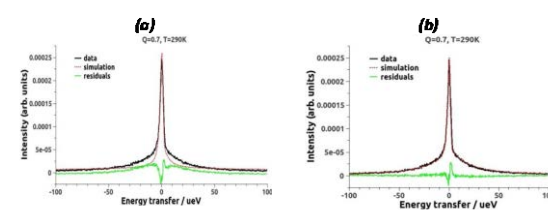


**Clearing house**

**EDEN**

**ParCAT**

**ESGF**

## An integrated compute and data environment for neutron science

- Provide live feedback from experiment through streaming data capture and reduction
- Integrating simulation (MD/DFT) directly into the neutron scattering data analysis chain



Example: , *ab-initio* MD simulations for ferroelectrics/thermoelectrics. Focus on *width* of dispersions

Fit between experiment (black) and simulated (red) dynamics structure factors for water-hydrogen.

# Case Study: Neutron Scattering
## Enabling near real-time feedback from experiment

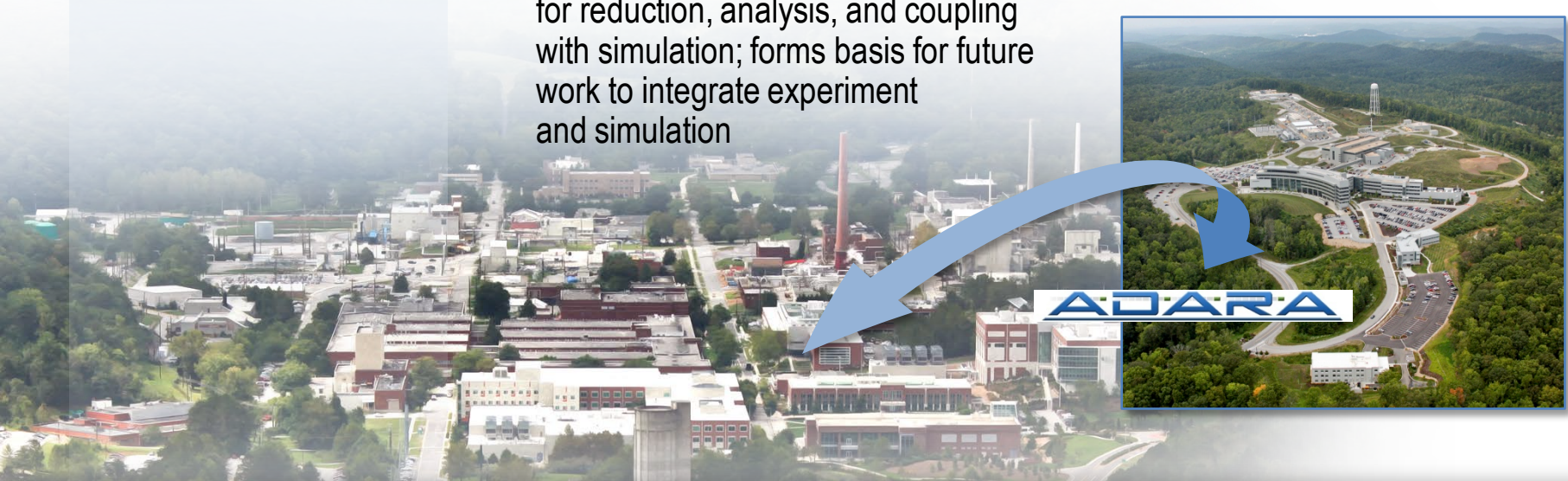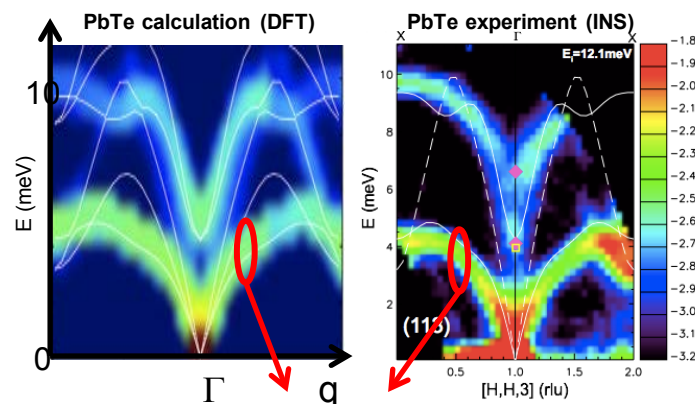| Challenge | Response | Status |
|---|---|---|
| Realizing the full potential of SNS requires near real-time feedback to users and integration of experiment and simulation/modeling | • Accelerating Data Acquisition, Reduction and Analysis (ADARA)<br><br>  – Leverages ORNL'S neutron scattering and computational expertise<br><br>  – Stream data to computational resources and provide live feedback from experiment in real-time<br><br>  – High performance data backplane for reduction, analysis, and coupling with simulation; forms basis for future work to integrate experiment and simulation | • ADARA is being field tested on the SNS HYSPEC beam line<br><br>  – Complete field testing on HYSPEC beam line<br><br>  – Continue development and deployment of ADARA on subsequent SNS beam lines |

**OAK RIDGE NATIONAL LABORATORY**
MANAGED BY UT-BATTELLE FOR THE U.S. DEPARTMENT OF ENERGY

# Neutron Science: Integrating Simulation In the Data Analysis Infrastructure

- The CAMM will integrate materials modeling/simulation (MD/DFT) directly into the chain for neutron scattering data analysis, **offline** and **online** (in near real time)

- Developing workflows for refinement, integration of MD codes, **neutron scattering corrections** ..

- The CAMM is working with ORNL's Materials Science and Technology Division to study coarse grained MD simulations of polymers PEO-AA (CNMS), *ab-initio* MD simulations for ferroelectrics/thermoelectrics



Example: *ab-initio* MD simulations for ferroelectrics/thermoelectrics. Focus on *width* of dispersions

---

**The Center for Accelerating Materials Modeling (CAMM)**

- *Partnership between ORNL's Neutron Sciences, Physical Sciences and Computing and Computational Sciences Directorates*
- *ORNL SEED money and DOE funds provided to study force field refinement from quasi-elastic and inelastic neutron scattering data*
- *CAMM formed in response to BES proposal call for* Predictive Theory and Modeling

# Case Study: Neutron Imaging
## Developing mathematical foundations for UQ-VV

APPLIED & COMPUTATIONAL

PMG

Predictive
Methods
Group

MATHEMATICS

**Mathematical Modeling**

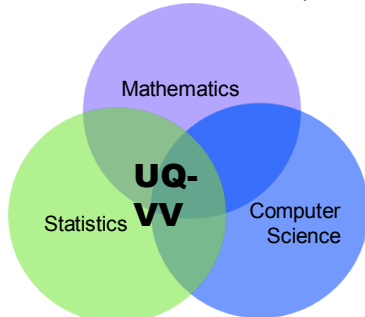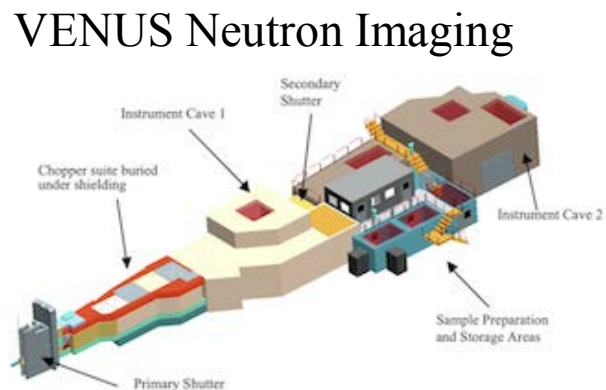**Statistical Challenge:** Data calibration and bias correction

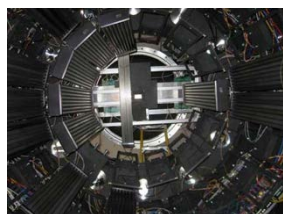**Math Challenge:** Effective model UQ representation

**Observation and Experiment**

**Computer Simulation**

VENUS Neutron Imaging

**Computation Science Challenge:** Architecturally effective embedded UQ-VV

ORNL Leadership Computing

Secondary Shutter

Instrument Cave 1

Chopper suite buried under shielding

Instrument Cave 2

Sample Preparation and Storage Areas

Primary Shutter

Mathematics

UQ-VV

Statistics

Computer Science

OAK RIDGE NATIONAL LABORATORY
MANAGED BY UT-BATTELLE FOR THE U.S. DEPARTMENT OF ENERGY

# Neutron Science: An Integrative Data Infrastructure

# Questions?

OAK RIDGE NATIONAL LABORATORY