# Data Management, Analysis and Visualization ASCAC Update

Lucy Nowell, PhD

Computer Scientist and Program Manager

Advanced Scientific Computing Research
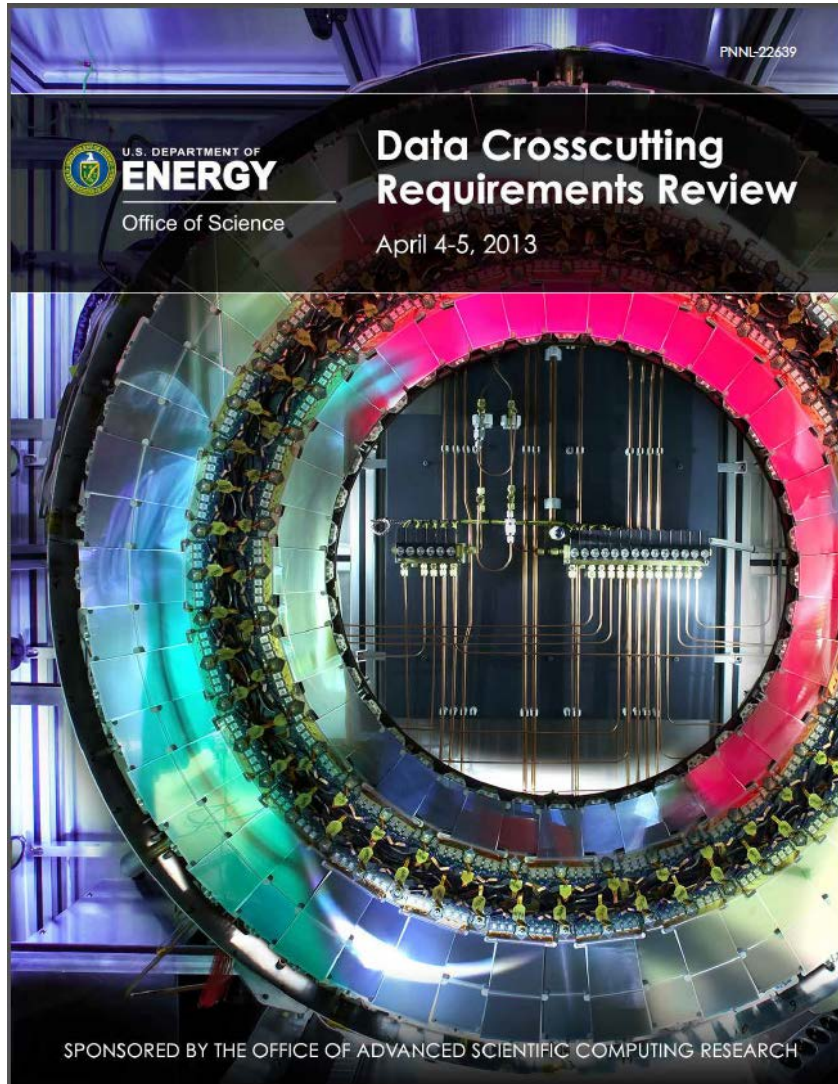
24 March 2015

# Overview

- **Summary of the 14-1053 Announcements on Scientific Data Management, Analysis and Visualization at Extreme Scale 2**

- **Portfolio Analysis and Awards Made**

- **Workshops**
  - Storage Systems and I/O (SSIO) Summit
  - Data Council Meeting
  - SSIO Workshop Series and forthcoming report

- **Data/Vis PI Meeting**

- **Next Steps**

- **Questions?**


- **Data Analytics and Visualization Sub-Plexus Leads:**
  - **Wes Bethel, LBNL, and Jim Ahrens, LANL**

- **Data Management Nexus Leads:**
  - **Rob Ross, ANL, and Gary Grider, LANL**

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# ASCR Computer Science Program

- **ASCR Core CS Program tries to address two fundamental questions:**
  - **How can we make today's and tomorrow's leading edge computers tools for science?**
  - **How do we extract scientific information from large data from experiments and simulation?**

- **There are several factors that provide important context for the ASCR Core CS program:**
  - ASCR Facilities (Leadership-class supercomputers at Argonne and Oak Ridge National Laboratories; capacity supercomputer (NERSC) at Lawrence Berkeley National Laboratory
  - Research and Evaluation Partnerships
  - ASCR's Applied Mathematics and Next Generation Networking Programs
  - SciDAC Centers and Institutes
  - Exascale Codesign Centers

U.S. DEPARTMENT OF **ENERGY** | Office of Science
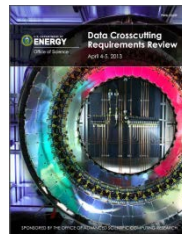
# ASCR and BES, BER, HEP



## Data Crosscutting Requirements Review

In April 2013, a diverse group of researchers from the U.S. Department of Energy (DOE) scientific community assembled in Germantown, Maryland to assess data requirements associated with DOE-sponsored scientific facilities and large-scale experiments.

http://science.energy.gov/~/media/ascr/pdf/program-documents/docs/ASCR_DataCrosscutting2_8_28_13.pdf

# Research Challenges

- How can data be represented in the system so as to maximize its analytic value while also minimizing the power and memory cost of the analytic process?

- How can data provenance, which is essential for validation and later reuse/repurposing, be captured and stored without overburdening a system that is input/output (I/O) bound?

- For complex scientific problems that require integrated analysis of data from multiple simulations, observatories, and/or disciplines, how can the expected IO and memory constraints be overcome to support re-use and repurposing of data?
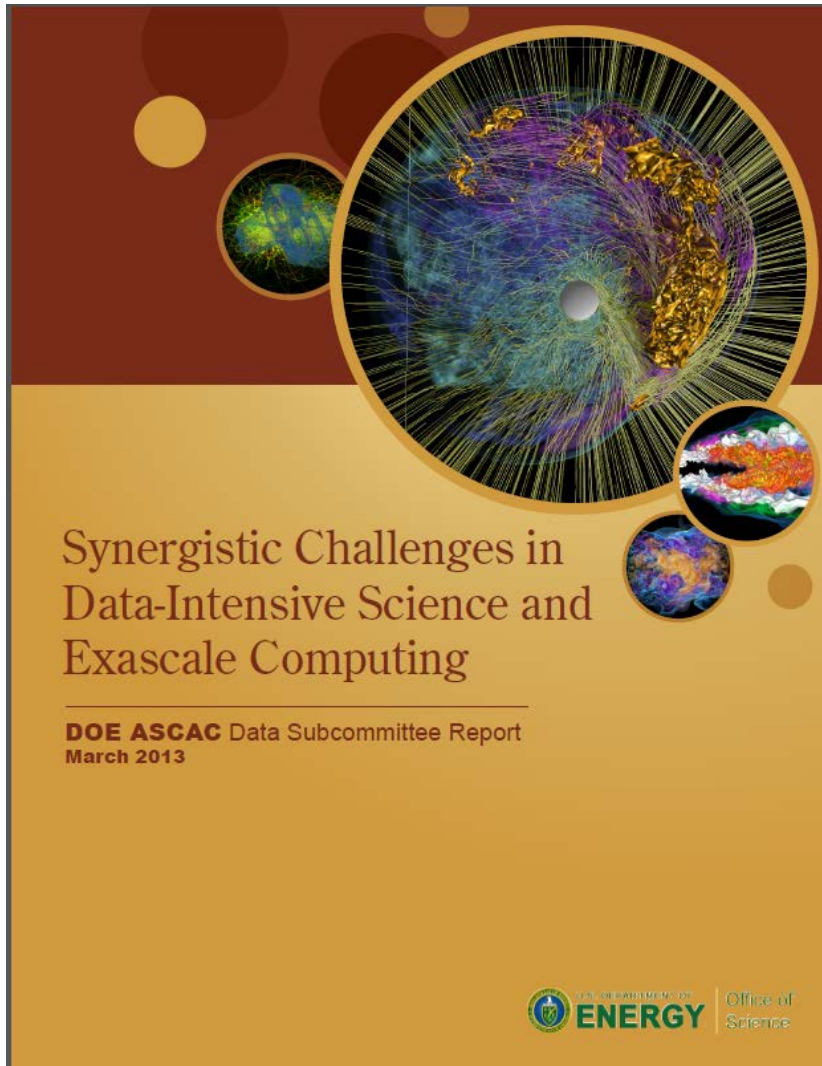
# Research Challenges (cont.)

- In the context of these memory and I/O constrained systems, how can simulation data be compared to or integrated with observational/experimental data, both to validate the simulations and to support new types of analysis?

- What new abstractions are needed for long-term data storage that move beyond the concept of files to more richly represent the scientific semantics of experiments, simulations, and data points?

- How can data analysis contribute to generating the ten to one hundred billion way concurrency that future machines will support and need to mask latency?

- How can data management and analysis applications help to mitigate the impact of frequent hardware failures and silent faults?

# DOE ASCR Advisory Committee (ASCAC)
## Data Subcommittee Report

Synergistic Challenges in Data-Intensive Science and Exascale Computing

**DOE ASCAC** Data Subcommittee Report
**March 2013**

This report discusses the natural synergies among the challenges facing data-intensive science and exascale computing, including the emergence of a new scientific workflow.

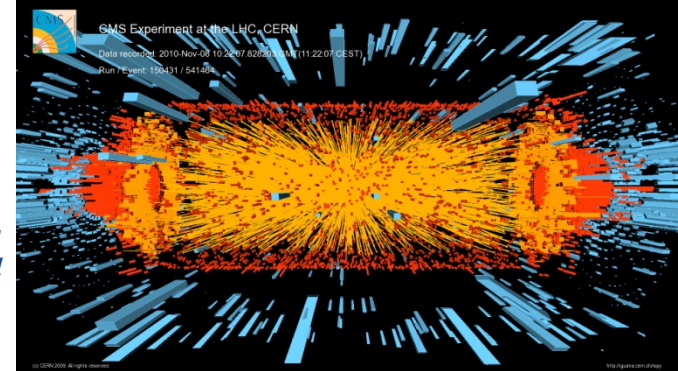http://science.energy.gov/~/media/ascr/ascac/pdf/reports/2013/ASCAC_Data_Intensive_Computing_report_final.pdf

# Exascale and Data-intensive Science
## *Interplay Between Data and Computing Critical to 21st Century Science*

- **DOE missions require ASCR to address both simultaneously**
  - Can't be decoupled even if we could afford to do so
- **Data-intensive science faces many of the same technology challenges of exascale**
  - Some are even worse for "big-data"
  - Energy use is the grand challenge (e.g. the square kilometer array estimates 100MW needed for computing)
- **ASCAC charge**
  - Subcommittee looking at DOE mission needs, big data and exascale to identify synergies and high priority research needs



*Petabytes per second*



*Petabytes today, Exabytes tomorrow*



*"… it would be a mistake to think of them ["big data" and "big compute"] as independent activities. Instead, their requirements are tightly intertwined since they both contribute to a shared goal of scientific discovery." (ASCAC report on Synergistic Challenges in Data-Intensive Science and Exascale Computing page 5)*

**U.S. DEPARTMENT OF ENERGY** | Office of Science

# *Scientific Discovery at the Exascale Workshop*



*Scientific Discovery at the Exascale: Report from the DOE ASCR 2011 Workshop on Exascale Data Management, Analysis and Visualization*, February 2011, Houston, TX

Organizer: Sean Ahern, ORNL; Co-Chairs: Arie Shoshani, LBNL, and Kwan-Liu Ma, UC Davis

http://science.energy.gov/~/media/ascr/pdf/program-documents/docs/Exascale-ASCR-Analysis.pdf

# Workshop Report

- **Principal Finding: "The disruptive changes posed by a progressive movement towards the exascale in HPC threaten to derail the scientific discovery process. Today's success in extracting knowledge from large HPC simulation output are not generally applicable to the exascale era, and simply scaling existing techniques to higher concurrency is not sufficient to meet the challenge." – p. 1**



**Scientific Discovery at the Exascale:**

Report from the DOE ASCR 2011 Workshop on Exascale Data Management, Analysis, and Visualization

February 2011
Houston, TX

Workshop Organizer:
  Sean Ahern, Oak Ridge National Laboratory
Co-Chairs:
  Arie Shoshani, Lawrence Berkeley National Laboratory
  Kwan-Liu Ma, University of California Davis
Working Group Leads:
  Alok Choudhary, Northwestern University
  Terence Critchlow, Pacific Northwest National Laboratory
  Scott Klasky, Oak Ridge National Laboratory
  Kwan-Liu Ma, University of California Davis
  Valerio Pascucci, University of Utah

Additional Authors:
Jim Ahrens          Kenneth Moreland
E. Wes Bethel       George Ostrouchov
Hank Childs         Michael Papka
Jian Huang          Venkatram Vishwanath
Ken Joy             Matthew Wolf
Quincey Koziol      Nicholas Wright
Gerald Lofstead     Kesheng Wu
Jeremy Meredith

Sponsored by the Office of Advanced Scientific Computing Research

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# 14-1043 Scientific Data Management, Analysis & Visualization at Extreme Scale 2

- **The purpose of this announcement was to invite proposals for basic computer science research on five major themes:**

  1. **Usability and user interface design;**

  2. *In situ* **methods for data management, analysis and visualization;**

  3. **Design of** *in situ* **workflows to support data management, processing, analysis and visualization;**

  4. **New approaches to scalable interactive visual analytic environments; and/or**

  5. **Proxy applications or workflows and/or simulations for data management, analysis and visualization software to support co-design of extreme scale systems**

- **The supported research will lay the foundation for building the software infrastructure to support scientific data management, analysis and visualization in the context of extreme scale computing**. While described in the context of extreme scale supercomputing and simulations, these topics also have **broad implications for data intensive science at scale and have been called out as synergistic challenges**.

# Additional Guidance to Applicants & Reviewers

- Proposals may address one or more of the announcement themes. Preference will be given to larger collaborative projects that offer more comprehensive and coherent research and solutions, provided that such proposals are appropriately integrated into a whole and not a collection of disjointed efforts.

- Proposals that address Theme 1 must show **clear relevance to challenges in scientific data management, analysis and visualization in the context of science disciplines supported by the DOE Office of Science and the ASCR supercomputing environment at ANL, ORNL and LBNL**.

- Proposals addressing Themes 2-5 **must explicitly address the challenges of operating within the expected exascale environment**, including severe constraints on data movement, worsening I/O bottleneck, frequent hard and soft faults, and the necessity for extremely high levels of concurrency.

- Proposals that address a combination of Theme 1 with one or more of Themes 2-5 are welcome.

- Proposals must clearly establish that the proposed research supports the missions of the Advanced Scientific Computing Research office and the Computer Science program, as well as showing relevance to other Office of Science disciplines. **Proposed research must also be relevant to future extreme scale computing platforms operated by ASCR and must advance pertinent aspects of computer science.**

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# Not Covered in This Solicitation

- **The 14-1043 did not address a number of related challenges:**
  - File systems and input/output (FSIO)
  - Distributed data management and analysis
  - Management and analysis of experimental and observational data
  - Use of data management, analysis and visualization for system administration, program development and debugging, and/or use by the underlying operating and runtime systems

- **We hope to issue another solicitation in FY 15 that will address these topics.**
  - Three projects from 10-256 in the area of FSIO were renewed in August of 2013

- **Also out of scope:**
  - Topics covered in Applied Mathematics solicitation on data at extreme scale.

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# Review Criteria

- Proposals will be subjected to scientific merit review (peer review) and will be evaluated against the following criteria, listed in descending order of importance.

  1. Scientific and/or Technical Merit of the Project;

  2. Appropriateness of the Proposed Method or Approach;

  3. Competency of Applicant's Personnel and Adequacy of Proposed Resources;

  4. Reasonableness and Appropriateness of the Proposed Budget; and

  5. Relevance to the mission of and systems operated by the Advanced Scientific Computing Research program

# Expected Award Sizes

- **The overall budget was approximately $4M per year for three years with potential to increase to $7+M per year, depending on proposal quality**

- Each funded project could have a total annual budget of $150,000 for a small project at a single institution up to $1.5 million for a project spanning a larger portion of the scope of research and multiple institutions.

- An award to non-lab applicant could not exceed $500K per year for three years if partnered with a DOE National Lab and could not exceed $350K per year if not partnered with a DOE National Lab.

# Award Selection

- **The Selection Officials considered the following items:**

  - **Scientific and technical merit of the proposed activity as determined by merit review**

  - **Availability of funds**

  - **Relevance of the proposed activity to Office of Science priorities**

  - **Ensuring an appropriate balance of activities within Office of Science and ASCR programs**

  - **Previous performance**

  - **Relevance to the mission of and systems operated by the Advanced Scientific Computing Research program.**

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# Summary Statistics

- **FOA/Lab Announcement 14-1043 posted 12/19/13**

- **Pre-proposals due 2/7/14:**
  - 74 unique pre-proposals received

- **Encourage/Discourage date 2/21/14:**
  - 39 projects/pre-proposals encouraged

- **Encouraged proposals due 4/2/14:**
  - 35 encouraged projects' proposals received and reviewed
  - Proposal review panel May 7-8, 2014

- **9 projects recommended for award**
  - 3 University led
  - 6 Lab led
  - 35 awards in all processed:
    - 16 DOE National Labs, 15 University, 4 Industry

- **Total funding awarded:**
  - $24.852M spanning three years
  - ~ 67% to Labs, 25% Universities; 8% Industry

# Funded

- *Usable Data Abstractions for Next-Generation Scientific Workflows*
  - PI: Deb Agarwal, LBNL, with ORNL, University of California Berkeley, and University of Washington

- *Optimizing the Energy Usage and Cognitive Value of Extreme Scale Data Analysis Approaches*
  - PI: Jim Ahrens, LANL, with Virginia Tech, University of New Hampshire, and University of Texas Austin (2 awards)

- *Scalable Analysis Methods and In Situ Infrastructure for Extreme Scale Knowledge Discovery*
  - PI: Wes Bethel, LBNL, with ANL, Georgia Tech, JMSI, Inc., and Kitware, Inc.

- *A Unified Data-Driven Approach for Programming In Situ Analysis and Visualization*
  - PI: Pat McCormick, LANL, with SNL, University of Utah, Stanford University, and Kitware, Inc.

# Funded Projects

- ***XVis: Visualization for the Extreme-Scale Scientific-Computation Ecosystem***
  - PI: Ken Moreland, Sandia, with ORNL, LANL, University of California Davis, University of Oregon, and Kitware, Inc.

- ***High-Performance Decoupling of Tightly Coupled Data Flows***
  - PI: Tom Peterka, ANL, and SNL

- ***In Situ Indexing and Query Processing of AMR Data***
  - PI: Nagiza Samatova, NCSU, with LBNL

- ***Performance Understanding and Analysis for Exascale Data Management Workflows***
  - PI: Karsten Schwan, Georgia Tech, with ORNL and University of Oregon

- ***Extreme-Scale Distribution-Based Data Analysis***
  - PI: Han-Wei Shen, Ohio State University, with ANL and LANL

**U.S. DEPARTMENT OF ENERGY** | Office of Science

# Usable Data Abstractions for
# Next-Generation Scientific Workflows
## PI: Deb Agarwal, LBNL

- **Co-PIs:**
  - S. Vazkudai, ORNL; M. Franklin, University of California Berkeley; C. Aragon, University of Washington

- **Research Themes:**
  - 1. Usability and User Interface Design
  - 3. In-Situ Data Management and Workflow Management Systems
  - 5: Proxy Applications and/or Simulation Support for Co-Design

- **Sciences: Material Sciences, Climate, Combustion Physics, Sustainable Systems**

- **Technical Summary:**
  - The project will focus on
    - Identifying data abstractions using empirical ethnographic studies,
    - Exploring two levels of data abstractions (user and system) that enable user provided suggestion to help manage efficiency trade-offs in storage and data-management, and
    - Developing representative mini-workflows for use in studying the performance and energy characteristics.

U.S. DEPARTMENT OF ENERGY | Office of Science

# Optimizing the Energy Usage and Cognitive Value of Extreme Scale Data Analysis Approaches
## Jim Ahrens, LANL

- **Co-PIs:**
  - David Rogers and Jon Woodring, LANL; Colin Ware, University of New Hampshire; Greg Abram, University of Texas Austin; Francesca Samsel, University of Texas Austin

- **Research Themes**
  - 1. Usability and user interface design and
  - 2. In situ methods for data management, analysis and visualization

- **Sciences: Oceanography, Cosmology, and Plasma Physics**

- **Technical Summary:**
  - This project is specifically focused on developing general workflow, data and user evaluations that will impact a broad range of scientific domains.
  - A goal of the project is to significantly improve the way analysis algorithms are developed, measured and used at the petascale and in the future at the exascale. The research aims to move beyond simple performance measurements to assess the value of analysis algorithms, ultimately achieving perceptual optimization of visualization algorithms and cognitive optimization of interactive methods in visualization tools. The project is specifically focused on developing general workflow, data and user evaluations that will impact a broad range of scientific domains, largely through adaptive sampling techniques and power profiles for those techniques on emerging testbeds and current DOE petascale production systems.

U.S. DEPARTMENT OF ENERGY | Office of Science

# Scalable Analysis Methods and In Situ Infrastructure for Extreme Scale Knowledge Discovery
## PI: Wes Bethel, LBNL

- **Co-PIs: G. Weber, LBNL; Vishwanath, ANL; Wolf, Georgia Tech; Duque, Intelligent Light; O'Leary, Kitware Inc.**

- **Research Theme: 2: In Situ Data Management, Analysis, and Visualization**

- **Technical Summary:**

  - This project will develop new algorithms for analysis, and visualization – topological, geometric, statistical analysis, flow field analysis, pattern detection and matching – suitable for use in an *in situ* context and aimed specifically at enabling scientific knowledge discovery in several exemplar application areas of importance to DOE.

  - Complementary to the *in situ* algorithmic work, the project will focus on several leading *in situ* infrastructures and investigate research questions germane to enabling new algorithms to run at scale across a diversity of existing in situ implementations.

  - The intent is to move the field of *in situ* processing towards portability, so that it may ultimately be possible to write an algorithm once and then have it execute in one of several different *in situ* software implementations.

# A Unified Data-Driven Approach for Programming *In Situ* Analysis and Visualization

## Pat McCormick, LANL

- **Co-PIs: J. Bennett, Sandia; B. Geveci, Kitware, Inc.; C. Hansen, University of Utah; A. Aiken, Stanford University**

- **Unfunded Advisors: J. Chen, Sandia; K. Heitmann, ANL; P. Jones, LANL; S. Keckler, NVIDIA; M. Schulte, AMD**

- **Research Themes:**

  - 2) In situ methods for data management, analysis and visualization;

  - 3) Design of in situ workflows to support data management, processing, analysis and visualization; and

  - 5) Proxy applications or workflows and/or simulations for data management

- **Technical Summary**

  - The overarching goal of the proposed effort is the development of an unified data-driven approach for programming applications and *in situ* analysis and visualization. The project will study the impact of:

    - Supporting effective *in situ* data management, analysis and visualization;

    - Providing a foundation for building efficient and effective workflow management;

    - Enabling an interactive *in situ* user environment on the underlying runtime software design; and

    - Understanding the impact of these options on existing applications, infrastructure and tools.

# XVis: Visualization for the Extreme-Scale Scientific-Computation Ecosystem
## Ken Moreland, Sandia

- **Co-PIs: B. Geveci, Kitware, Inc.; J. Meredith, ORNL; C. Sewell, LANL; K.M. Ma, University of California Davis; H. Childs, University of Oregon**

- **Research Themes:**
  - 1) Usability and user interface design,
  - 2) In situ methods for data management, analysis and visualization,
  - 4) New approaches to scalable interactive visual analytic environments and
  - 5) Proxy applications

- **Technical Summary:**
  - The team will
    - Build a framework that enables research and development of massively threaded visualization algorithms and apply this framework to extreme-scale visualization;
    - Address the emerging challenges of in situ visualization including both resource costs and practical application;
    - Assess the usability of extreme-scale visualization techniques by conducting user studies that measure the understanding visualization techniques impart;
    - Build proxy applications and use these to empirically study the complicated interactions of visualization with simulation codes and heterogeneous architectures, and
    - Use this insight to drive their algorithmic and in situ research.

# High-Performance Decoupling of Tightly Coupled Data Flows
## Tom Peterka, ANL

- **Co-PI: Jay Lofstead, Sandia**

- **Research Themes**
  - 2) *In situ* methods for data management, analysis and visualization;
  - 3) Design of *in situ* workflows to support data management, processing, analysis and visualization, and
  - 5) Proxy applications or workflows and/or simulations for data management, analysis and visualization software to support co-design

- **Technical Summary:**
  - This proposal seeks to reduce the dependence of existing workflow systems on such one-off data translation functions by developing standard interfaces with general and usable functionality. The researchers have designed Decaf to be a small library of standardized reusable dataflow tools sitting below workflows in the software stack so that Decaf will assist workflows in their work. This project will augment the team's dataflow with essential data operators— pipelining, selection, and aggregation—in order to achieve space-time data permutations within the dataflow, and they will execute the dataflow over various transport layers in current and future HPC architectures. The project will deliver three things:
    - A library of dataflow primitives, akin to "BLAS for dataflow,"
    - A method for automatically constructing broadly applicable dataflows from the same set of primitives, and
    - A generic and reusable solution that other workflow and coupling tools can use.

# *In Situ* Indexing and Query Processing of AMR Data
## Nagiza Samatova, NCSU

- **Co-PI: S. Byna, LBNL**
- **Research Theme: 2) *In situ* Data Management, Analysis and Visualization**
- **Science driver: Climate modeling**
- **Technical Summary:**
  - The proposed *in situ* query framework for AMR data is a multi-faceted approach that entails:
    - AMR-aware, *in situ* indexing and meta-indexing to handle the unique adaptive data structure;
    - Interactive and *in situ* query processing over AMR data to yield near-real-time scientific insights and simulation steering feedback; and
    - An AMR-tailored query model and API to bind the aforementioned services to end-user scientific applications.
  - Additionally, the research will explicitly address extreme-scale system requirements in the development of these concepts in order to ensure the continued applicability of this approach. The proposal asserts that this will be the first effort to provide such a systematic in situ AMR query framework.

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# Performance Understanding and Analysis
# for Exascale Data Management Workflows
## Karsten Schwan, Ga. Tech

- **Co-PIs:  A. Maloney, University of Oregon; H. Abbasi, ORNL**

- **Research Themes:**
  - 2) *In Situ* Data Management, Analysis and Visualization, and
  - 5) Proxy Applications and/or Simulation Support for Co-Design

- **Technical Summary:**
  - The project will explore the actual performance characteristics of exemplar in situ workflows for applications that include Particle-in-Cell (PIC), Particle, and Finite Difference (FD) kernels.
  - The goal is to determine how the resulting in situ performance workflows can be integrated with these data intensive application workflows and how dynamic tuning of data accesses impacts end-to-end workflow performance. Further, the project seeks to determine how this performance data can be used to create models for the different actors in these workflows, so that the researchers can then generate end-to-end workflow proxies that will accurately emulate the behavior of fully composed applications. The generated proxies will be one of the artifacts of this research and will enable more in-depth collaboration with hardware vendors, OS researchers, and the DOE co-design centers.

# Extreme-Scale Distribution-Based Data Analysis
## Han-Wei Shen, Ohio State University

- **Co-PI/Senior Personnel: T. Peterka, ANL; J. Woodring & J. Wendelberger, LANL; G. Agrawal & H. Wang, Ohio State University**
- **Research Theme: 2)** *In situ* **methods for data management, analysis and visualization**
- **Sciences: Climate, Cosmology, and Superconductivity.**
- **Technical Summary:**
    - The key development of the project will be a distribution-based analysis and visualization framework based on in situ processing of extreme-scale scientific data. Project goals are to ensure that scientists can easily obtain an overview of the entire data set regardless of the size of the simulation; understand the characteristics and locations of the features; easily interact with the data and select regions and features of interest; and perform all the analysis tasks with a small memory footprint. To achieve these goals, the analysis framework will consist of three unique but tightly integrated components:
        - Computation, representation, and indexing of data distributions generated from large-scale simulation data;
        - Spatial, value, and temporal domain data summarization, reduction, and triage with distributions; and
        - Efficient exploratory data analysis and visualization based on distributions.

# Storage Systems and I/O Summit

- **Day long event held Sept. 15, 2014, in Rockville, MD**
- **Organized by Rob Ross of ANL and Gary Grider of LANL**
- **Participants: 8 DOE Lab representatives, including Rob & Gary**
- **Focus on identifying ways to reinvigorate research in this area**
- **Outcome:**
  - Planning for the SSIO Workshops held in December 2014

# Data Council Meeting

- **Two day event held Sept. 16-17, 2014 in Rockville, MD**

- **Organized by Jim Ahrens of LANL, with help from Wes Bethel of LBNL**

- **Participants: 20 DOE Lab representatives, spanning multiple areas of Computer Science and Applied Mathematics**

- **Focus was on the Exascale Preliminary Planning documents on Data Management and Data Analysis & Visualization**

- **Outcomes:**
  - Community concurrence with the EPPDD documents
  - Steps towards defining the Data Stack for exascale
  - Recommendations on coordination across the portfolio

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# DAV Exascale Software Stack

**Science Applications: Climate, Accelerator, Fusion, Materials, Biology, …**

**Scientific Workflow**

**Science-facing DAV applications and tools: VisIt, ParaView, R, Matlab, …**

I/O libraries
Simulation I/O
EOD ingest

**DAV Algorithms: Statistical Analysis, Machine Learning, Scientific and Information Visualization, Graph Analytics, Data Mining, …**

**DAV Libraries: VTK-m, DAX, EAVL, PISTON, R, …**

**Big Data Motif/Kernels: Structured memory, unstructured memory, graph methods, dense/sparse linear algebra, optimization, …**

**Optimized Libraries/Runtime: MapReduce, ScaLAPACK, cuBLAS, TBB/Thrust, MPI, …**

**Platform & Hardware: Many-/multi-core chipsets, deep memory hierarchy, power management, reliability management, storage, data movement**
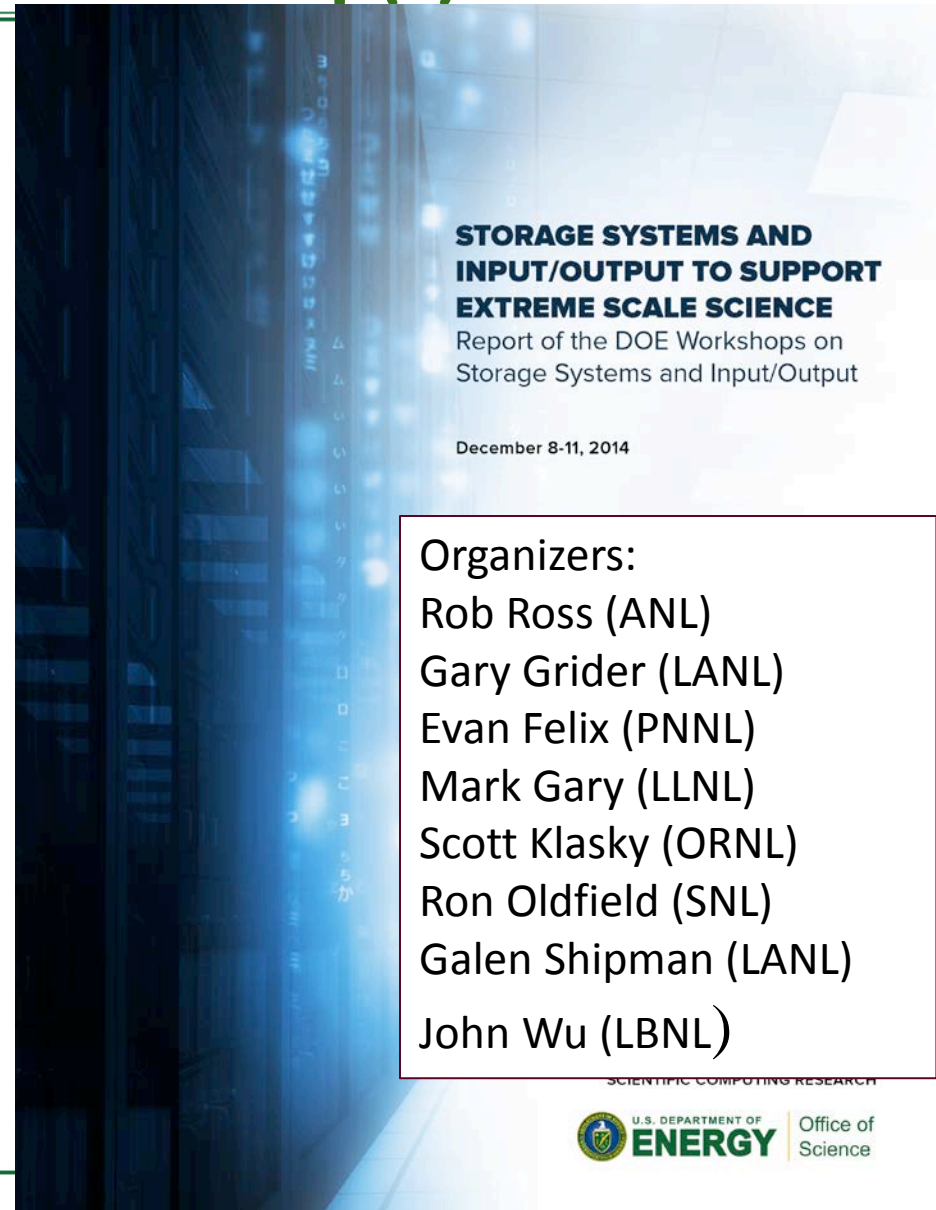
# Storage Systems and I/O Workshops

- **Organized by Rob Ross of ANL and Gary Grider of LANL, with support from SSIO Summit participants**

- **Four events held Dec. 8-11, 2014 in Rockville, MD**

  - **Science Requirements Review on Dec. 8**
    - Focused on scientific use cases for SSIO and patterns of data movement

  - **Crosscutting CS Review on Dec. 9**
    - Focus on identifying dependencies across areas and coordination needed to see that needs are met

  - **1.5 day Storage Systems and I/O Workshop on Dec. 10-11**
    - Included briefings on science use cases/requirements, patterns of data movement, and expected future supercomputer architectures
    - Charged to assess the state of the art, research needed to address identified requirements, and approaches to reinvigorating the field
    - Report forthcoming; summary follows

  - **.5 day Burst Buffers Workshop on Dec. 11**
    - Focused on DOE requirements and how to address them in future vendor RFPs
    - Ongoing teleconferences and email discussion

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# Recent SSIO Workshop(s)

- **Workshop(s) back in December 2014**

  – Day 1: apps people talking with SSIO and software stack people

  – Day 2: software stack and SSIO people using apps requirements to talk about cross cuts

  – Days 3-4: SSIO researchers talk about needed R&D given previous 2 days input

**STORAGE SYSTEMS AND INPUT/OUTPUT TO SUPPORT EXTREME SCALE SCIENCE**
Report of the DOE Workshops on Storage Systems and Input/Output

December 8-11, 2014

Organizers:
Rob Ross (ANL)
Gary Grider (LANL)
Evan Felix (PNNL)
Mark Gary (LLNL)
Scott Klasky (ORNL)
Ron Oldfield (SNL)
Galen Shipman (LANL)

John Wu (LBNL)

SCIENTIFIC COMPUTING RESEARCH

U.S. DEPARTMENT OF ENERGY | Office of Science

U.S. DEPARTMENT OF ENERGY | Office of Science

# Topics in HW/SW Architectures for Storage and I/O

- **Networks**
  - Topology
  - Understanding Impact
- **Deep Storage Hierarchies**
  - Technology Integration
- **Nonvolatile**
  - Use: Burst Buffers, On-Node
  - Memory/Storage Convergence (incl. Prog. Models)
- **Active Storage**
  - Programmable FS/Storage
- **Resilience/Integrity/Availability**
  - End-to-end Data Integrity
  - Resilience to Component Failures
  - Availability
  - Coding Techniques
  - Checkpointing Strategies

- **Autonomics**
  - Garbage Collection
  - Adapting to User Demands/Behavior
- **Understandability (also another session)**
  - Architecting for Understandability
  - Acting on Prediction
- **Security**
  - Security and HPC Services
  - Avoiding Information Leakage
- **New (to us) Paradigms/Capabilities**
  - Key-Value
  - Scheduling of I/O
  - Coupling with Experimental/Observational Facilities

# SSIO Workshop Findings

- *In situ* data analysis is already an important component of many applications.
- New solid state and disk storage layers are complicating the storage hierarchy.
- Scientists need a coherent view and management methods of the storage resources.
- Current SSIO designs are hindered by their isolation from system-level resource management, monitoring, and workflow systems.
- Many important aspects of app/system SSIO behavior aren't well understood.
- New requirements for results validation may change the role of SSIO systems.
- New programming models/systems drive new persistence mechanisms.
- Scientists desire increasingly complex data abstractions that improve productivity.
- Community access to data on apps and systems needed, as well as test environments, for new technology evaluation and bringing new talent into the community.

# Near Term SSIO Research Priorities from Workshop

- **SSIO architectures research for:**
  - managing deep and heterogeneous storage hierarchies
  - alternative management paradigms to the file system model

- **In metadata, name spaces, and provenance:**
  - new methods of management of rich metadata
  - breaking away from the current file model

- **In the area of supporting science data :**
  - develop the next generation of I/O middleware and services to support new programming abstractions and workflows

- **In the area of understanding SSIO:**
  - improve our ability to characterize storage activities to model and predict the behavior of SSIO activities on future systems.

# Somewhat Bigger Picture: Data Management Services

| Users | Application Tasks | Analysis Tasks |
|---|---|---|

**Task and Data Coordination**

| Workflow Desc. and Mgmt. | Publish/ Subscribe |
|---|---|

| Science Data Model Services* | Pass-through |
|---|---|

| Programming Model |
|---|
| Resource Mgmt and Scheduling |
| Identity and Security |
| Performance Monitoring |
| WAN Data Services |

| Provenance Management | **Core Data Services** |
|---|---|

| Metadata Management | Core Data Model Services* |
|---|---|

| In System Storage | External Storage | Networking HW |
|---|---|---|

# Functionality in an SSIO solution

- **Data/metadata storage, persistence, resilience, naming**
- **Bridging between networks**
- **Performance monitoring**
- **Executing user code (presumably in some sandbox)**
- **Probably built using a custom RPC system**
- **Possibly built using a custom programming language**

- **Block management on storage resources**
- **Group membership, STOMITH (kill misbehaving resource)**
- **Client interface (usually embedded in node OS)**

- **Lots of potential for sharing of components with runtime/OS**

U.S. DEPARTMENT OF **ENERGY** | Office of Science    **Slide by Rob Ross, ANL**    Nowell – ASCAC SDMAV Briefing 24 March 2014    38

38

# Some Questions (in Conclusion)

- **In what ways are runtime systems consumers of storage and I/O resources/services?**
  - What data models and capabilities are needed to support those uses?
- **How do runtimes and SSIO cooperate to connect the pieces of the deep memory hierarchy?**
- **How can SSIO, runtimes, and other components capture adequate provenance for validation of results?**
- **What's the relationship between runtimes and active storage?**
- **What common building blocks are shared by runtimes and SSIO? Can technology be shared?**
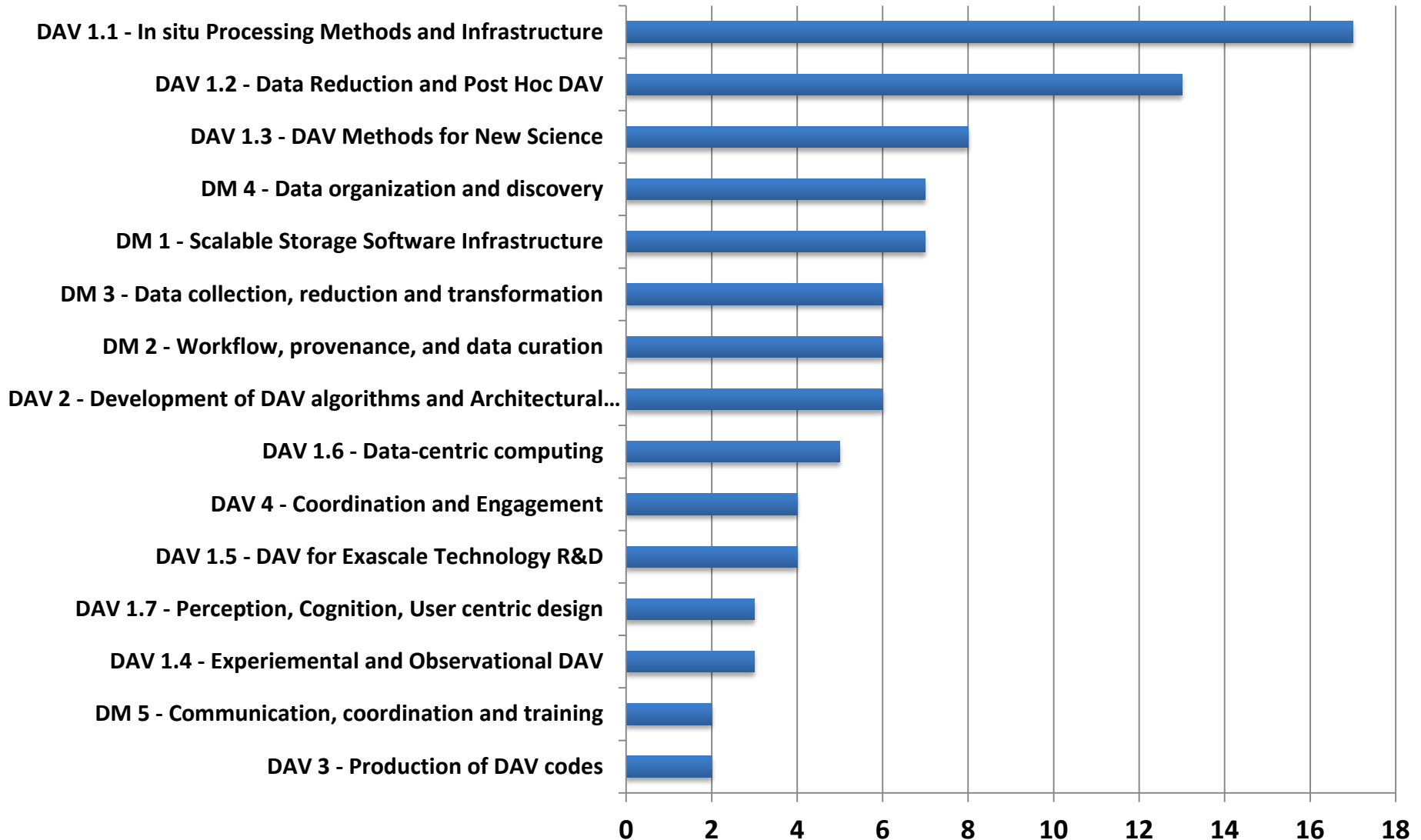- **How do we coordinate cross-cutting activities?**

U.S. DEPARTMENT OF **ENERGY** | Office of Science
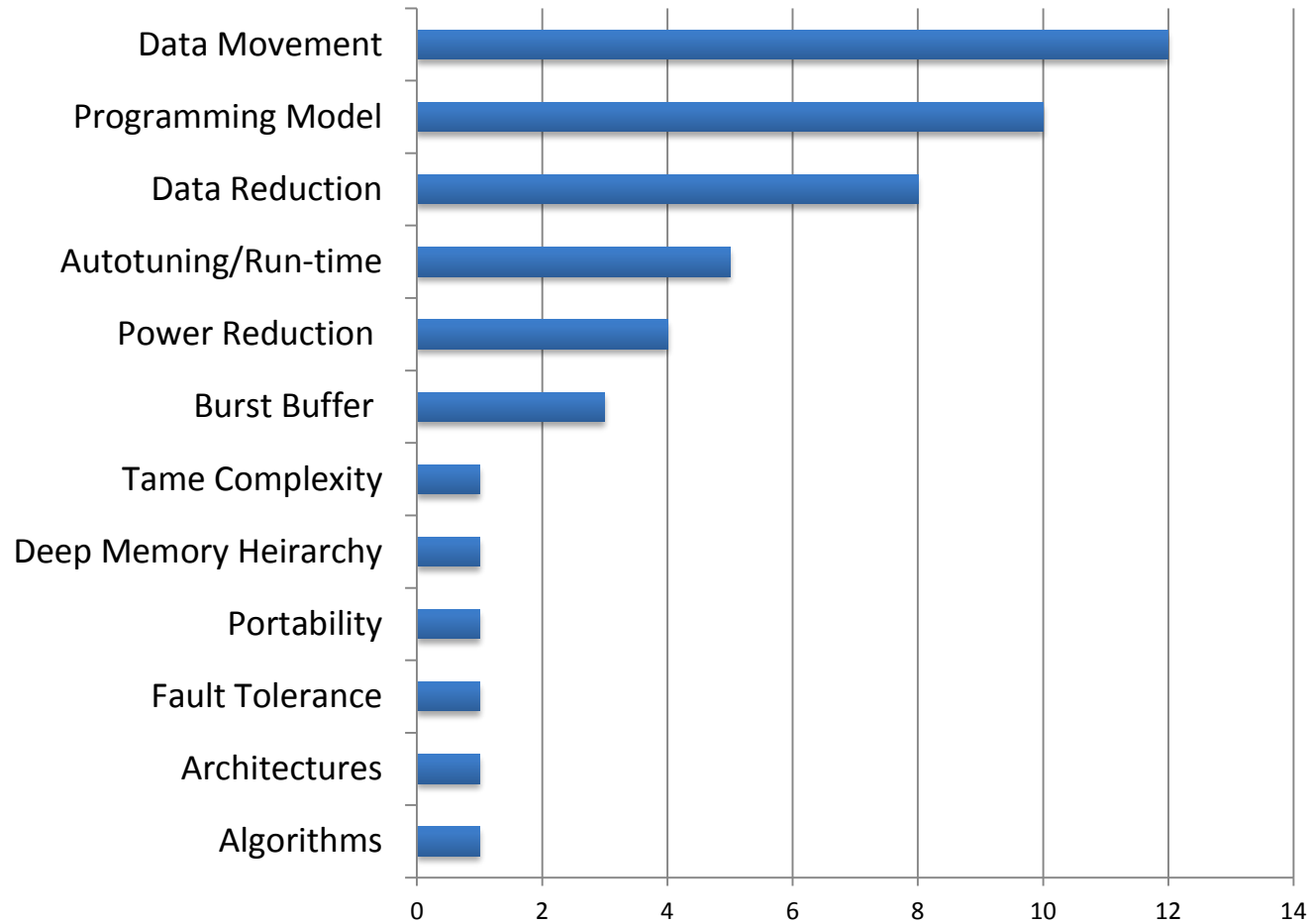
# Data/Vis PI Meeting

- **Organized by Wes Bethel of LBNL, with support from Jim Ahrens of LANL**

- **Three day event held Jan. 13-15, 2015, in Walnut Creek, CA**

- **http://extremescaleresearch.labworks.org/events/2015-jan-scientific-data-management-analysis-and-visualization-pi-meeting**

- **~80 participants spanning ASCR CS DM and DAV portfolio**

- **Goals and Objectives:**

  - Nexus/Plexus meetings, coordination and community collaboration

  - Examine the EPPDD sections on DM and DAV and assess how funded research fits within that framework

  - Understand and evaluate impact of future supercomputer architectures with respect to the research projects

  - Discuss the Data Stacks and map research projects into them

  - Begin discussion of the integration and further development that will be needed to support scientific discovery at extreme scale

U.S. DEPARTMENT OF ENERGY | Office of Science

# Exercise 1: Map projects into the Exascale Plan



Bar chart showing project counts:

- DAV 1.1 - In situ Processing Methods and Infrastructure: 17
- DAV 1.2 - Data Reduction and Post Hoc DAV: 13
- DAV 1.3 - DAV Methods for New Science: 8
- DM 4 - Data organization and discovery: 7
- DM 1 - Scalable Storage Software Infrastructure: 7
- DM 3 - Data collection, reduction and transformation: 6
- DM 2 - Workflow, provenance, and data curation: 6
- DAV 2 - Development of DAV algorithms and Architectural...: 6
- DAV 1.6 - Data-centric computing: 5
- DAV 4 - Coordination and Engagement: 4
- DAV 1.5 - DAV for Exascale Technology R&D: 4
- DAV 1.7 - Perception, Cognition, User centric design: 3
- DAV 1.4 - Experiemental and Observational DAV: 3
- DM 5 - Communication, coordination and training: 2
- DAV 3 - Production of DAV codes: 2

U.S. DEPARTMENT OF ENERGY | Office of Science

# What technologies do projects provide for exascale?



Horizontal bar chart showing counts by technology:
- Data Movement: 12
- Programming Model: 10
- Data Reduction: 8
- Autotuning/Run-time: 5
- Power Reduction: 4
- Burst Buffer: 3
- Tame Complexity: 1
- Deep Memory Heirarchy: 1
- Portability: 1
- Fault Tolerance: 1
- Architectures: 1
- Algorithms: 1

# What do projects need from the community for success?

U.S. DEPARTMENT OF **ENERGY** | Office of Science

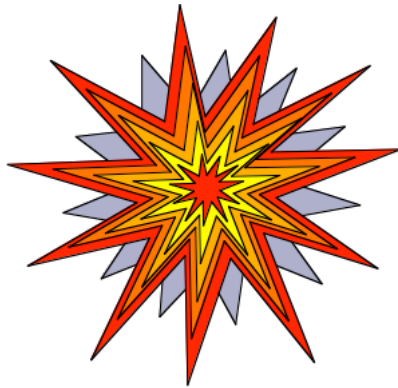# Exascale architectural impacts and challenges/opportunities identified from talks:

# Next Steps

- **Planning for Announcements in SSIO**

- **Coordination to ensure that SDMAV dependencies are met**
  - Nexus/Plexus leads attended March 2015 Software Stack workshops on Programming Models/Environments and Runtime Systems

- **Workshop on Workflows April 20-21, 2015, with Rich Carlson**
  - Organizers are Tom Peterka of ANL and Ewa Deelman of ISI

- **Workshops on requirements for experimental and observational science - TBD**

- **Planning for 2016 SDMAV PI Meeting**

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# Questions?

# Thank you!

# Lucy Nowell

lucy.nowell@science.doe.gov