Exec Committee

Pete Beckman

Jack Dongarra

Sergi Girona

Yutaka Ishikawa

Satoshi Matsuoka

Dan Reed

# White House HPC Initiative

**The White House**
Office of the Press Secretary

For Immediate Release                                July 29, 2015

## Executive Order -- Creating a National Strategic Computing Initiative

EXECUTIVE ORDER

- - - - - - -

CREATING A NATIONAL STRATEGIC COMPUTING INITIATIVE

By the authority vested in me as President by the Constitution and the laws of the United States of America, and to maximize benefits of high-performance computing (HPC) research, development, and deployment, it is hereby ordered as follows:

Section 1.  Policy.  In order to maximize the benefits of HPC for economic competitiveness and scientific discovery, the United States Government must create a coordinated Federal strategy in HPC research, development, and deployment.  Investment in HPC has contributed substantially to national economic prosperity and rapidly accelerated scientific discovery.  Creating and deploying technology at the leading edge is vital to advancing my Administration's priorities and spurring innovation.  Accordingly,

# NSCI has 5 Strategic Themes

Create systems that can apply exaflops of computing power to exabytes of data.

Keep the United States at the forefront of HPC capabilities.

Improve HPC application developer productivity

Make HPC readily available

Establish hardware technology for future HPC systems.

# International Exascale Software Project Meetings

**Overall goal:**

- Bring together the international community to explore plans and identify gaps for producing a software infrastructure capable of supporting exascale applications

**Meeting history:**

1. Santa Fe, NM, US, April 2009
2. Paris, France, June 2009
3. Tsukuba, Japan, October 2009
4. Oxford, UK, April 2010
5. Maui, HI, US, October 2010
6. San Francisco, US, April 2011
7. Cologne, Germany, October 2011
8. Kobe, Japan, April 2012

SC08 (Austin), SC09 (Portland), ISC09, SC10 (NOLA), ISC10, SC11 (Seattle)

www.exascale.org: White Papers and Slides

# IESP Roadmap (2009 – 2012)

The IESP Roadmap presented a multidimensional analysis of the major challenges to be overcome in order to create a software infrastructure capable of supporting exaflop performance on next generation systems, and made a cogent case for the urgency of starting that work as soon as possible.

Spurred in some degree by the work of the IESP and its Roadmap, the United States, the Europe an Union, and Japan have, in the past three years, moved aggressively to develop their own plans for achieving exascale computing in the next decade

# IESP Roadmap Components

www.exascale.org

# IESP ➡ BDEC

## BDEC derived much of its impetus from the earlier work

- International Exascale Software Project (IESP)
- European Exascale Software Initiative (EESI)
- European Exascale Software Initiative 2 (EESI2)
- European eXtreme Data and Computing Initiative (EXDCI)

# Big Data and Extreme Computing workshops series (BDEC)

http://www.exascale.org/bdec/

**Overarching goal:**

1.Create an international collaborative process focused on the **co-design of software infrastructure** for extreme scale science, addressing the **challenges of both extreme scale computing and big data**, and supporting a broad spectrum of major research domains,

2.Describe funding structures and strategies of public bodies with Exascale R&D goals worldwide

3.Establishing and maintaining a global network of expertise and funding bodies in the area of Exascale computing

BDEC Workshop, Charleston, SC, USA, April 29-May 1, 2013

BDEC Workshop, Fukuoka, Japan, February 26-28, 2014

BDEC Workshop, Barcelona, Spain, January 28-30, 2015

BDEC Workshop being Planning in the USA for March 2016

## 1 - BDEC Workshop, Charleston, SC, USA, April 29-May 1, 2013

Big Data and Extreme-scale Computing (BDEC) Workshop, Charleston, SC, USA, April 29-May 1, 2013

- Workflow Issues
- Architecture Challenges
- Higher Level Data Challenges : Data provenance, Policy based data management, Environments that support new types of data-driven research, Shared software infrastructure for intermediate processing
- Software Challenges: Tools to support real-time monitoring and observation of workflows, Coordination between data movement and compute services,  Mechanisms to support fault tolerant workflows in data analysis, Mini-apps to support infrastructure co-design, Integration of widely used BD-capable data libraries into standard packages, Common tools for managing and exploring data, Interoperability Challenges

---

BDEC Workshop Report (November 29, 2013)

### Report on the
### Big Data and Extreme-scale Computing (BDEC) Workshop,
### Charleston, SC, USA, April 29-May1, 2013

**1   Introduction**

This report on the Big Data and Extreme-scale Computing (BDEC) workshop offers an initial account of the effort to develop a plan for sustained international cooperation in the design and development of a new generation software infrastructure for extreme scale science. The meeting, the first of a planned series, derived much of its imperus from the earlier work of the International Exascale Software Project (IESP) and the European Exascale Software Initiative (EESI, http://www.eesi-project.eu). The goal of the IESP was two-fold: 1) to produce a plan for a common, high quality computational environment for the peta/exascale systems that are expected to arrive over the next decade; and 2) to mobilize and coordinate the work of the international open source software community to create that environment. BDEC retains those goals but changes the point of view. The EESI coordinated the European contribution to IESP.

The IESP, working through a series of eight international meetings held from 2009 to 2012, built on a range of important earlier studies, including [1-4], to produce a widely read and cited "roadmap" document. The IESP Roadmap [5] presented a multidimensional analysis of the major challenges to be overcome in order to create a software infrastructure capable of supporting exaflop performance on next generation systems, and made a cogent case for the urgency of starting that work as soon as possible. Spurred in some degree by the work of the IESP and its Roadmap, the United States, the European Union, and Japan have, in the past three years, moved aggressively to develop their own plans for achieving exascale computing in the next decade. The EESI produced a European roadmap along with a set of recommendations to address the Petascal/Exascale challenge [10].

The first BDEC workshop marks a beginning of a distinct new phase of this community movement. The motivation for this second stage is based on the recognition that the "digital data deluge," which was sighted on the horizon well over a decade ago [6], has finally made landfall with impressive force. It is apparent that in the era of "Big Data," when every major field of science and engineering is producing, and needs to (repeatedly) process, truly extraordinary amounts of data, the many unsolved problems surrounding *wide-area, multistage workflows—the diverse patterns of when, where, and how all that data is to be produced, transformed, shared, and analyzed*—have to take center stage. Although the IESP Roadmap shows a clear awareness that extreme scale science inevitably means extreme scale data as well as extreme scale computing, IESP working groups, for the most part, adopted a traditional HPC (i.e., supercomputer centric) perspective. They were largely (and understandably) preoccupied by the impending software crisis caused by the move to the new paradigm in hardware and systems architecture, a paradigm that demands orders of magnitude more parallelism, places unprecedented constraints on energy consumption, and requires resilience to faults occurring at far higher frequencies than ever before. Data-driven workflow issues received some collateral discussion in the Roadmap, but the focus of attention for the IESP was on the revolutionary innovations in the system software stack that would be needed to address the steep challenges of emerging peta/exascale systems.

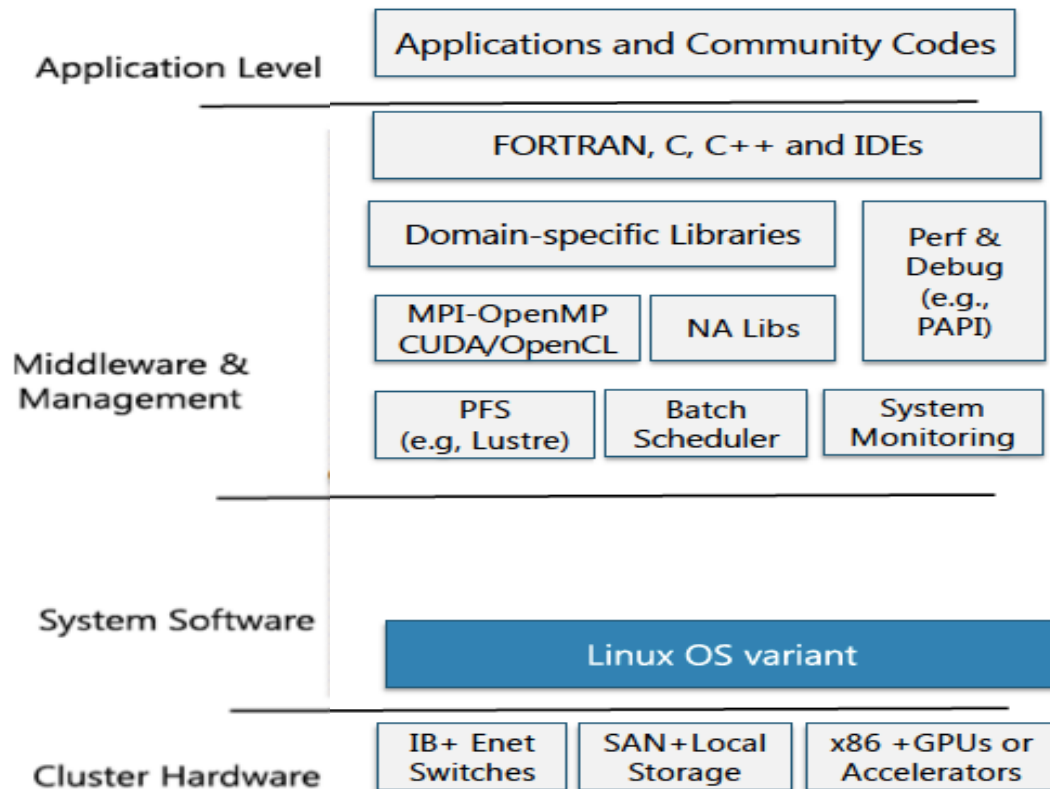# Big Data and Extreme Computing

High-end data analytics and HPC are both essential elements of an integrated computing research-and-development agenda;

- Big compute generates and is needed to analyze big data
- Networking and memory performance are critical to both

Programming models and tools are perhaps the biggest point of divergence between the scientific-computing and big-data ecosystems.

# EcoSystems



| Application Level | Applications and Community Codes |
|---|---|

FORTRAN, C, C++ and IDEs

Domain-specific Libraries | Perf & Debug (e.g., PAPI)

Middleware & Management

MPI-OpenMP CUDA/OpenCL | NA Libs

PFS (e.g, Lustre) | Batch Scheduler | System Monitoring

System Software | Linux OS variant

Cluster Hardware | IB+ Enet Switches | SAN+Local Storage | x86 +GPUs or Accelerators

Computational Science Ecosystem

As scientific research increasingly depends on both high-speed computing and data analytics, the potential interoperability and scaling convergence of these two ecosystems is crucial to the future.



| | Data Analytics Ecosystem | Computational Science Ecosystem |
|---|---|---|
| **Application Level** | Mahout, R and Applications | Applications and Community Codes |
| **Middleware & Management** | Zookeeper (coordination), Hive, Pig, Sqoop, Flume, Map-Reduce, Storm, Hbase BigTable (key-value store), HDFS (Hadoop File System), AVRO, Cloud Services (e.g., AWS) | FORTRAN, C, C++ and IDEs; Domain-specific Libraries; Perf & Debug (e.g., PAPI); MPI-OpenMP CUDA/OpenCL; NA Libs; PFS (e.g, Lustre); Batch Scheduler; System Monitoring |
| **System Software** | VMs, Containers and Cloud Services; Linux OS variant | Linux OS variant |
| **Cluster Hardware** | Ethernet Switches; Local Node Storage; Commodity X86 Racks | IB+ Enet Switches; SAN+Local Storage; x86 +GPUs or Accelerators |

## Computational Science

**FORTRAN,C,C++**: languages
**PAPI**: performance and debugging tool
**MPI/OpenMP**: multi-core parallel model
**SLURM**: batch scheduler
**Lustre**: parallel file system

## Data Analytic

**Mahout**: machine learning tool
**Hive**: data warehouse software
**Pig**:  provide high level language for big data
**Sqoop**: exchange data with traditional database
**Flume**: log management
**Zookeeper**: maintaining consistency
**Storm**: real-time computation system.
**Hbase**: a distributed, scalable big data store.
**AVRO**:  data serialization system.

**NOTE:  The Divergence of Big Data and HPC Eco-Systems!**

# How Did We Get Here

Previous BDEC meetings: US & JP

Application Drivers: Astronomy, Medical, Genomics, Climate, Human Brain, Satellite images (GIS), Social Networks, etc.

Good discussions on converged / shared problems:

- Architecture, operations, software stack, algorithms

The BDEC effort *carries forward the general* mission of the IESP the co - design of software infrastructure for extreme scale science drawing on international cooperation and supporting a broad spectrum of major research domains; but it reframes the problems involved to fully take account of varied (and evolving) workflow patterns that such different communities of inquiry must create to work with both data and computing resources that are unprecedented in their scale

# Broad Goals

Develop ways for EC and BD communities to work closely together

Develop a *shared* vision for the future

- Architecture & Facilities
- Software
- Applications

Identify key research areas and develop a roadmap for building and extending BDEC capabilities in support of science