



# Improving HPC Software

Pete Beckman & Jack Dongarra

# Outline



- Current State: HPC Software
- Background: Activities in Europe and Japan
- The Changing Architecture
- The IESP Workshops
- Roadmap and Outcomes

# The Open Source Community Provides Most of the World's HPC Software



Jaguar	Total	XT5
Peak Performance	1,645	1,382
AMD Opteron Cores	181,504	150,176
System Memory (TB)	362	300
Disk Bandwidth (GB/s)	284	240
Disk Space (TB)	10,750	10,000
Interconnect Bandwidth (TB/s)	532	374



## National Energy Research Scientific Computing Center (NERSC)

- Located at Lawrence Berkeley National Lab
  - Cray XT-4 Franklin: 102 Tflop/s, 9,660 nodes, 19,320 cores
  - IBM Power 5 Bassi: 6.7 Tflop/s, 888 cores
  - Linux Opteron Cluster Jacquard: 3.1 Tflop/s, 712 cores
- Franklin quad-core currently being upgraded 350 Tflop/s, 38,640 cores
- NERSC-6 Project
  - RFP issued in September 2008
  - Proposals are being reviewed



Franklin



Bassi



Argonne's IBM Blue Gene/P – 556 TFs

# The Community is Diverse and Robust

- Over the last 10 years, the galvanization of the Open Source movement has dramatically improved HPC software

A very small sample:

- Linux Operating System, libc
- Python, Perl
- PAPI, TAU, Kojak
- UPC
- MPICH, OpenMPI
- ScaLAPACK
- VisIt
- GASNet, ARMCI/GA
- PVFS
- CFEngine, bconfig
- Ganglia
- SLURM, Cobalt
- Dyninst
- Torque/Moab, OpenPBS
- Charm++
- pNetCDF, HDF5
- GridFTP
- FFTW

# A Long History of Collaboration



The Result....

# Open Source HPC Software Stacks for Small Linux Clusters are Everywhere

Innovating@Sun Community Voices How to Buy Log In

United States [Change] English



Downloads & Trials Products Services Solutions Support Training Sun For...

Search

Home > Products > Software > Enterprise Computing > HPC and Grid Computing > Sun HPC Software >

## Sun HPC Software, Linux Edition

### Open Rich and Verified

Download Sun HPC Software, Linux Edition for free

Overview Features Tech Specs Support **Get It**

#### Get Sun HPC Software, Linux Edition

Download Sun HPC Software, Linux Edition today at no cost.

#### Free Sun HPC Software, Linux Edition 1.2 Download

Sun HPC Software, Linux Edition 1.2 is available to download now. Please note that technical support is not included with the software.

#### What You Get

- Lustre 1.6.6
- perfctr 2.6.36
- Env-switcher 1.0.13
- genders 1.9
- git 1.6.0.4
- Heartbeat 2.1.4-2.1
- Mellanox Firmware tools 2.5.0
- Modules 3.2.6
- MVAPICH 1.0.1
- MVAPICH2 1.0.3
- OFED 1.3.1
- OpenMPI 1.2.6
- RRDTool 1.2.26
- OpenSM 3.0.3
- pdsh 2.16
- Powerman 1.0.32
- HPCC Bench Suite 1.2.0
- Lustre IOKit
- IOR 2.10.1
- LNET self test
- NetPIPE 3.7.1
- Slurm 1.3.10
- MUNGE 0.5.8
- Ganglia 3.0.7
- oneSIS 2.0.1
- Cobbler 1.0.3
- CFEngine 2.2.6
- Conman 0.2.1
- FreeIPMI 0.6.6
- IPMTool 1.8.9.1
- lshw B.02.12.01

#### Clustertech HPC Environment Software Stack

Comprehensive, robust and scalable operation environment for every Linux HPC cluster.

##### Features and Benefits

- **Centralized cluster administration:** Cluster management made easy with centralized web and command line interfaces.
- **Real-time system monitoring:** Robust monitoring infrastructure and alert systems.
- **Rich set of components:** Covering all the components you need, including resource manager, portal middleware and tools for distributed access.
- **Seamless integration:** All peripheral components are tightly integrated with the core infrastructure and need not be individually configured.
- **Scalable, superior scalability:** by the FlexConf role-based cluster configuration engine.
- **High compatibility:** Supports 32 and 64-bit architectures. Supports leading interconnect technologies.
- **Hardware fault tolerant (HFTS):** Professional only. Eliminates the risk of single point failure by advanced high availability technologies.

CHES increases the efficiency, improves the price/performance ratio, reduces the TCO and maximizes the returns of your HPC investments by providing a comprehensive, robust and truly scalable environment for your HPC cluster.

##### CHES Overview

CHES addresses the three major software aspects of HPC clusters:

- The configuration of the operating system (networks, services, access control, etc.) at every node to form a coherent HPC platform.
- The convenient monitoring, control and management in the daily operation of the HPC cluster.
- The provision of tools and libraries for development and execution of HPC applications.

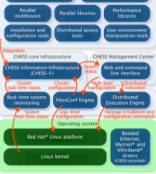
The central layer of CHES consists of two components: the CHES Information Infrastructure (CHES-ii) and the CHES management center. CHES-ii handles low-level system configuration issues, as well as consolidated real-time system status for monitoring purposes. Used in conjunction with the Management Center, CHES-ii allows an effective administration of the HPC cluster from an intuitive yet flexible interface.

CHES provides in its top layer all commonly used tools and libraries for the development and execution of HPC applications. These components are tightly integrated with the core infrastructure; they need not be individually reconfigured after system configuration changes, which significantly reduce maintenance effort.

##### The Real Scalability and Maintainability

Many cluster solutions on the market are advertised to be scalable, extensible and having low maintenance costs. However, the following technical difficulties of scaling and maintaining an HPC cluster are often overlooked:

- Adding more execution nodes to a cluster may not yield the expected result as bottlenecks that degrade the scalability of an application could be created by the application. For instance, the gateway for connection to



clustertech.com |

## Platform™ The Power of Sharing

my.platform.com hpcommunity.org

Products Industries Services Partners Resources Company Newsroom

Home / Products / Platform Open Cluster Stack 5

### Platform Open Cluster Stack 5

- › Intel Cluster Ready
- › Additional Components
- › Workload Management
- › Features and Benefits
- › Software and Supported Hardware
- › Services and Support
- › White Papers
- › FAQ



### Platform Open Cluster Stack (OCS) 5

Minimize the cost and time spent on deploying and managing a Linux cluster

"Intel and Platform Computing have optimized the performance of hundreds of thousands of enterprise High Performance Computing nodes."

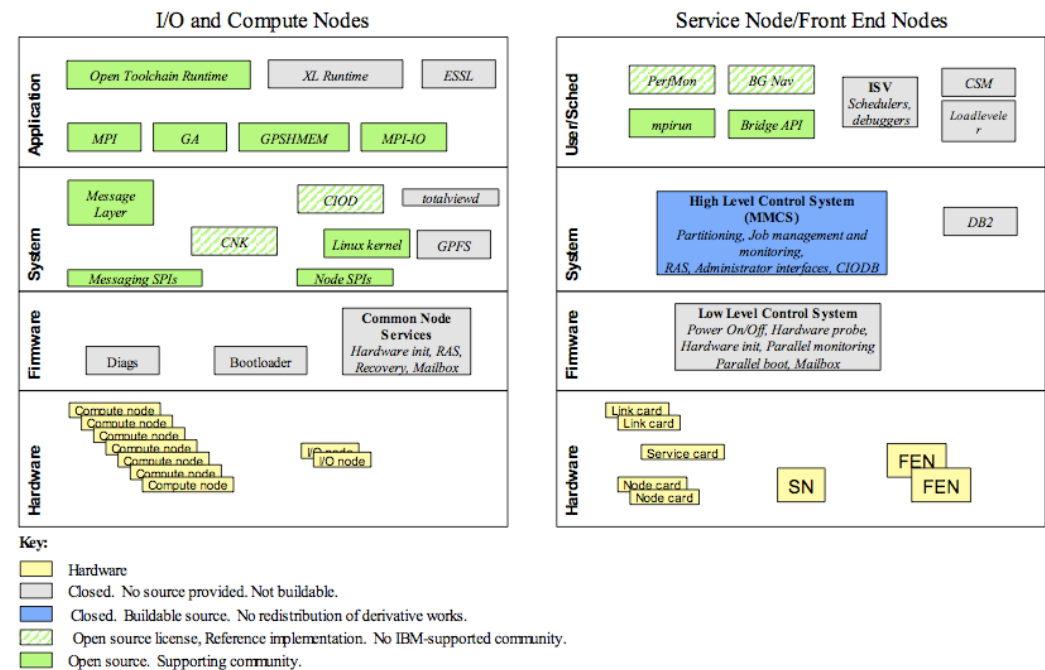
Richard Dracott, General Manager, Intel High Performance Systems



# Got Scale?

- For some markets, a closed source business model continues to work well
  - ▣ Single-node optimized math libraries & compilers
  - ▣ Debuggers for small clusters
  - ▣ Some queuing systems, parallel file systems, HSMs
  - ▣ Small cluster applications: Fluent, CFD++, etc

## BG/P Software Stack Source Availability



# Why Seek to Improve This?



- The largest scale systems are becoming more complex, with designs supported by large consortium
  - ▣ The software community has responded slowly
- Significant architectural changes arriving
  - ▣ Software must dramatically change
- Our ad hoc community coordinates poorly, both with other software components and with the vendors
  - ▣ Computational science could achieve more with improved development and coordination





# Mar 2012... the Japanese “> 10 PF” SC to be Online



**25~30MW Power**  
**~1mil ft2 floorspace**  
**\$1 bil construction**

**Kobe, Japan**  
**“the site”**

Interview with Ryutaro Himeno, Dr. Eng.  
Development Group Director at the Next-Generation  
Supercomputer R&D Center

A 10-petaflop\*  
supercomputer will be built  
in 2012. The  
supercomputer will boost a  
computing speed that is 50  
times faster than the  
world's current fastest  
computer. RIKEN is  
responsible for the  
development, construction,  
and operation of the  
supercomputer, which is a

huge government project with a total budget of 110 billion yen including facilities for the system and the research grid project. We interviewed Dr. Ryutaro Himeno, Development Group Director at the newly-established Next-Generation Supercomputer R&D Center to ask about why a speed of 10 petaflops, and what they are hoping to achieve with a 10-petaflop supercomputer.



*HIMENO Ryutaro*



Home

About PRACE

Activities

Use cases

Documents

Press corner

HPC Training

Contact us

### PRACE newsletter

Your e-mail address

HTML

Text

Subscribe

## Welcome to PRACE

The Partnership for Advanced Computing in Europe prepares the creation of a persistent pan-European HPC service, consisting of several tier-0 centres providing European researchers with access to capability computers and forming the top level of the European HPC ecosystem. PRACE is a project funded in part by the EU's 7th Framework Programme.

Supercomputers are indispensable tools for solving the most challenging and complex scientific and technological problems through simulations. To remain internationally competitive, European scientists and engineers must be provided with leadership-class supercomputer systems. PRACE, the Partnership for Advanced Computing in Europe will create a persistent pan-European high performance computing (HPC) service and infrastructure. This infrastructure will be managed as a single European entity. European scientists and technologists will be provided world-class leadership supercomputers with capabilities equal to or better than those available in the USA and Japan. The service will comprise three to five superior HPC centers strengthened by regional and national supercomputing centers working in tight



### News

- » HPC Infrastructures for Petascale Applications – DEISA PRACE Symposium 2009 2009-02-25
- » PRACE hosts highly successful Winter School 2009-02-25
- » PRACE held All Hands Meeting in Jülich, Germany, February 12-13, 2009 2009-02-16
- » Serbia joins the PRACE initiative 2009-02-12
- » PRACE Part of Zero-In Magazine – Call for Papers Open 2009-02-03 [more...](#)

### Events

- » OGF25 / EGEE User Forum, 2-6 March, Catania, Italy
- » 24th Forum ORAP, 26 March, Lille, France

# Traditional Sources of Performance Improvement are Flat-Lining (2004)

- New Constraints
  - 15 years of *exponential* clock rate growth has ended
- Moore's Law reinterpreted:
  - How do we use all of those transistors to keep performance increasing at historical rates?
  - Industry Response: parallelism doubles every 18 months *instead* of clock frequency!

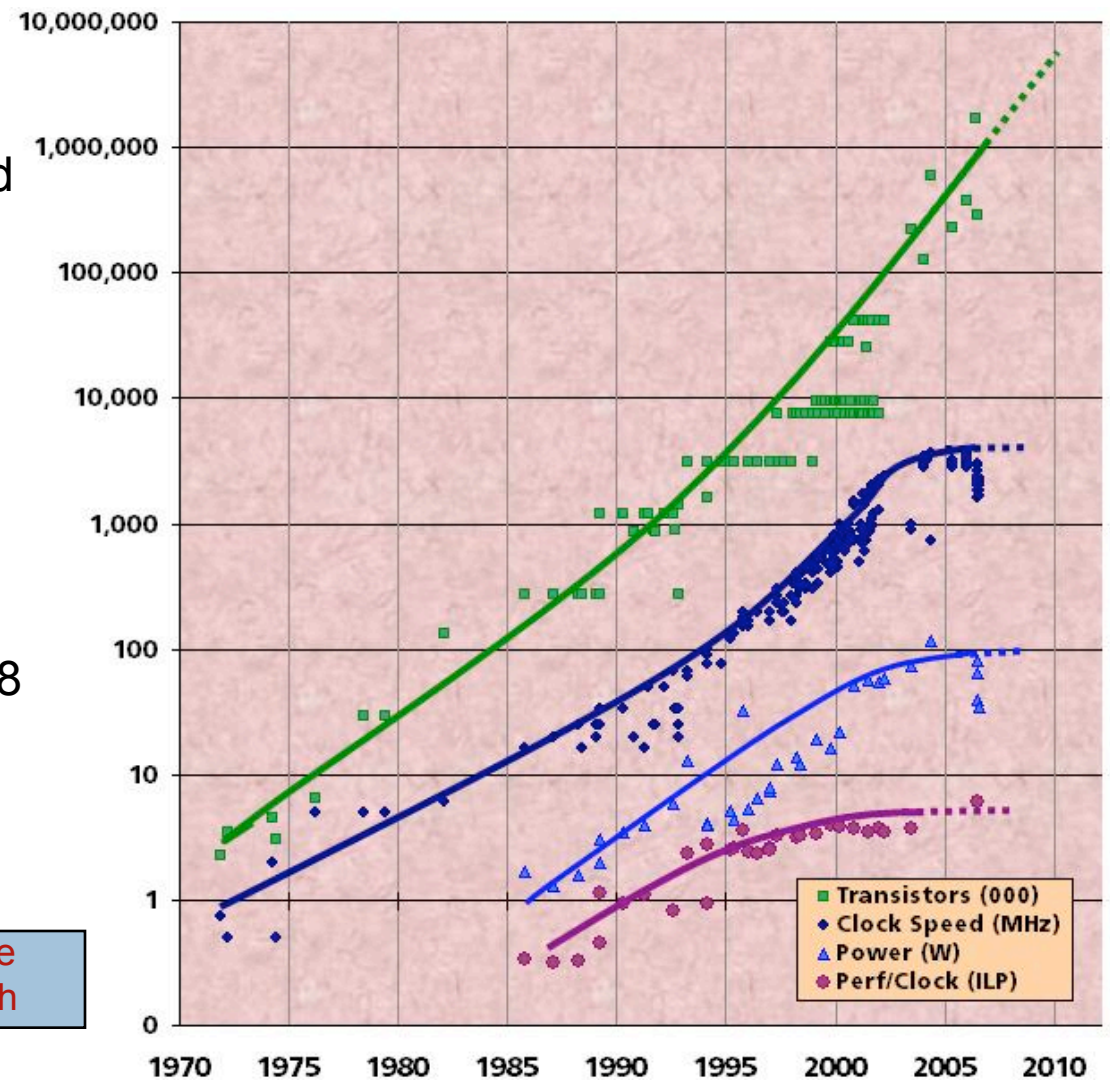
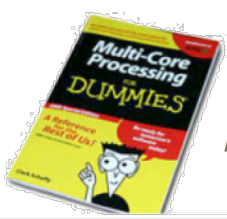
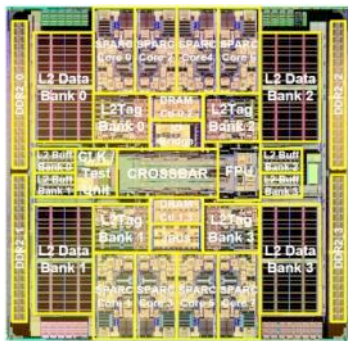


Figure courtesy of Kunle Olukotun, Lance Hammond, Herb Sutter, and Burton Smith



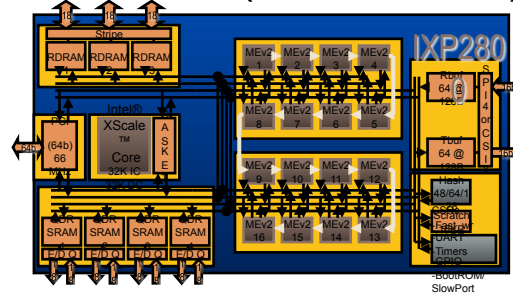
# Multicore comes in a wide variety

- Multiple parallel general-purpose processors (GPPs)
- Multiple application-specific processors (ASPs)

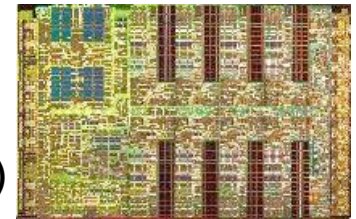
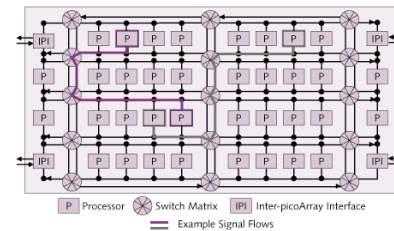


Sun Niagara  
8 GPP cores (32 threads)

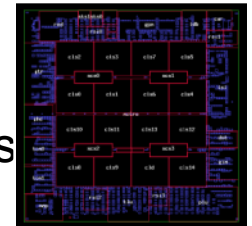
Intel Network Processor  
1 GPP Core  
16 ASPs (128 threads)



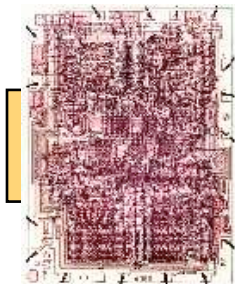
IBM Cell  
1 GPP (2 threads)  
8 ASPs



Picochip DSP  
1 GPP core  
248 ASPs



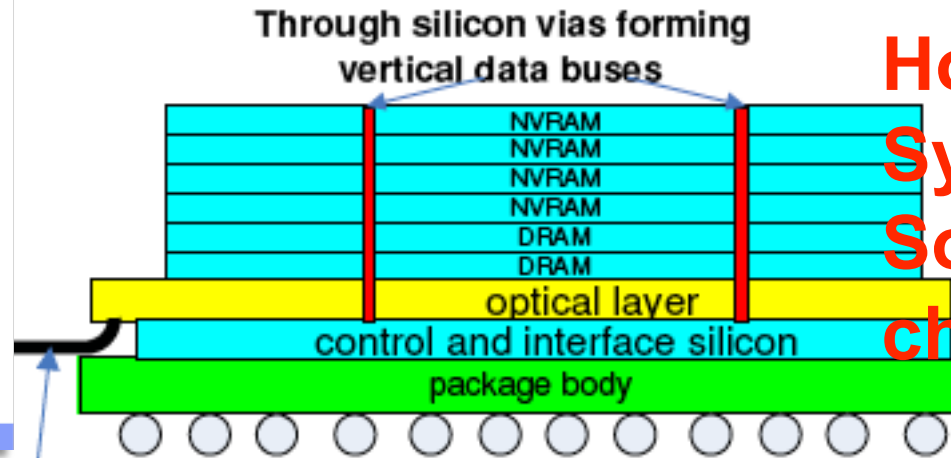
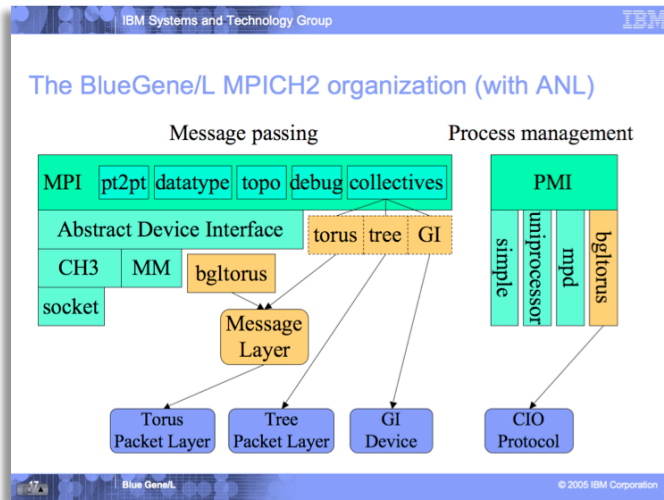
Cisco CRS-1  
188 Tensilica GPPs



Intel 4004 (1971):  
4-bit processor,  
2312 transistors,  
~100 KIPS,  
10 micron PMOS,  
11 mm<sup>2</sup> chip

*"The Processor is the new Transistor" [Rowen]*

# 3D Packaging: Changing Paradigms



**How will System Software change?**

Approach	Comments
<b>Distributed 3D stacks</b> <p>Direction of Heat Extraction</p> <p>CPU: 200 Cores (1/4 of total)</p> <p>4 DRAM die</p> <p>Interconnect Substrate</p>	<p>Distribute CPU across multiple memory stacks</p> <p>Assumes sufficient inter-stack bandwidth can be provided in substrate</p> <p>Likely to detract from performance, depending on degree of memory scatter</p>
<b>Advanced 3DIC</b> <p>CPU: 750 Cores</p> <p>16-32 DRAM die in groups</p> <p>Interposers</p> <p>Vias</p> <p>Substrate</p>	<p>Incorporate interposers into a single 17-33 chip stack to help in power/ground distribution and heat removal.</p> <p>Assumes Through Silicon Vias for signal I/O throughout chip stack</p>

Figure 7.5: Potential directions for 3D packaging (A).

Approach	Comments
<b>Advanced 3D Package</b> <p>CPU: 750 Cores</p> <p>16-32 DRAM die in groups</p> <p>Interposers</p> <p>Substrate</p> <p>Dense Via Field</p>	<p>To avoid complexity of a 33-chip stack, this approach, uses the interposers for high density signal redistribution, as well as assisting in power/ground distribution and heat removal.</p> <p>Requires a planar routing density greater than currently provided in thin film carriers.</p>
<b>Tiled Die</b> <p>CPU</p> <p>MEM</p> <p>CPU</p> <p>MEM</p>	<p>Use proximity connection or Through Silicon Vias to create memory bandwidth through overlapping surfaces.</p> <p>OR</p> <p>Tile with high bandwidth edge interfaces, using quilt packaging or an added top metal process. (Note, impact on latency and I/O power).</p>

Figure 7.6: Potential directions for 3D packaging (B).

# Where We Are Today:

We are not prepared for the changes coming

- Hardware features are uncoordinated with software development
  - ▣ (power mgmt, multicore tools, math libraries, advanced memory models, etc)
- Only basic acceptance test software is delivered with platform
  - ▣ UPC, HPCToolkit, Optimized libraries, PAPI, can be YEARS late
- Vendors often “snapshot” key Open Source components and then deliver a stale code branch
  - ▣ Counterexample: A model that works – MPICH for BG/P
- Community codes unprepared for sea change in architectures
- Coordination via SOW/contract is poor and only involves 2 parties
- No global evaluation of key missing components

# The IESP Workshops:



- Goal: Improve the world's simulation and modeling capability by improving the coordination and development of the HPC software environment.
  - ▣ Build a plan for how the international community can join together to improve software available for high-end systems over the next 2 to 10 years.
- The DOE, NSF, and EU have committed their support for the workshops.
- The first workshop will be Santa Fe, April 7-8.
  - ▣ White papers encouraged



# International Community Effort

- We believe this needs to be international collaboration for various reasons including:
  - The scale of investment
  - The need for international input on requirements
  - Europeans, Asians, and others are working on their own software that should be part of a larger vision for HPC.
- The process must be totally open

## **Executive Committee:**

Co-Chair: Jack Dongarra, Univ, of Tennessee / ORNL, US

Co-Chair: Pete Beckman, Argonne National Laboratory, US

Franck Cappello, INRIA, FR

Thomas Lippert, Jülich Supercomputing Centre, DE

Satoshi Matsuoka, Tokyo Institute of Technology, JP

Paul Messina, Argonne National Laboratory, US

# An Example Development Community

## The Apache Software Foundation *Meritocracy in Action.*



The Apache Software Foundation provides support for the Apache community of open-source software projects. The **Apache projects** are characterized by a collaborative, consensus based development process, an open and pragmatic software license, and a desire to create high quality software that leads the way in its field.

**We consider ourselves not simply a group of projects sharing a server, but rather a *community of developers and users.***

This page will give you everything you always wanted to know about the foundation but were afraid to ask. The difference between membership and committership, who decides what, how elections take place, how is our infrastructure setup, what is the board, what is a PMC, what's the philosophy behind the incubator, why is the foundation moving away from project containment. Come and see behind the scenes of the ASF.

- o **What is the Apache Software Foundation?**
- o **A bit of history**
- o **Meritocracy**
- o **The Foundation structure**
- o **Roles**
- o **Project management**
- o **The Foundation Infrastructure**
- o **The Foundation Incubator**
- o **Other Foundation entities**
- o **Conclusions**

### **What is the Apache Software Foundation?**

The Apache Software Foundation (ASF) is a 501(c)3 non-profit organization incorporated in the United States of America and was formed primarily to:

- o provide a foundation for open, collaborative software development projects by supplying hardware, communication, and business infrastructure
- o create an independent legal entity to which companies and individuals can donate resources and be assured that those resources will be used for the public benefit
- o provide a means for individual volunteers to be sheltered from legal suits directed at the

### **Apache Projects**

- o **HTTP Server**
- o **Abdera**
- o **ActiveMQ**
- o **Ant**
- o **APR**
- o **Archiva**
- o **Beehive**
- o **Camel**
- o **Cayenne**
- o **Cocoon**
- o **Commons**
- o **Continuum**
- o **CouchDB**
- o **CXF**
- o **DB**
- o **Directory**
- o **Excalibur**
- o **Felix**
- o **Forrest**
- o **Geronimo**
- o **Gump**
- o **Hadoop**
- o **Harmony**
- o **HiveMind**
- o **HttpComponents**
- o **iBATIS**

### **Foundation**

- o **FAQ**
- o **Licenses**
- o **News**
- o **Public Record:**
- o **Sponsorship**
- o **Donations**
- o **Thanks**
- o **Contact**

### **Foundation Projects**

- o **Conferences**
- o **Infrastructure**
- o **JCP**
- o **Legal Affairs**
- o **Security**
- o **Travel Assistance**

### **How it works**

- o **Introduction**
- o **Meritocracy**
- o **Structure**
- o **Roles**
- o **Collaboration**
- o **Infrastructure**
- o **Incubator**
- o **Other entities**

# Apache Foundation



- Create a foundation for open, collaborative software development projects by supplying hardware, communication, and business infrastructure
- Incubator projects can become Apache projects
- 800 “committers”
- The ASF Infrastructure is mostly composed of the following services:
  - ▣ the web serving environment (web sites and wikis)
  - ▣ the code repositories
  - ▣ the mail management environment
  - ▣ the issue/ bug tracking
  - ▣ the distribution mirroring system

# A Plan Could Include:



- Work with vendors to create the HPC equivalent to the ITRS (Int'l Tech Roadmap for Semiconductors)
  - Get community working on software before machine becomes available
- Community proposed unified roadmap for exascale software
- Identify missing components for future architectures and a plan to address them
- Develop models for working more closely with vendors
  - ▣ (support, acceptance tests, target features)
- Identify key application areas to drive development
- Community software development models
- Funding and organizational models (Apache, etc)

# Achievable Outcomes



- Improve the capability of computational science
- Build and strengthen international collaborations and leadership; deliver more capable, productive HPC systems
- Build and improve R&D program developing new programming models and tools addressing extreme scale
- Open source HPC development guided by roadmap with better coordination and fewer missing components
- Joint programs in education and training for the next generation of computational scientists.
- Vendor engagement and coordination for more capable software supporting exascale science

# Workshops and Report



- 3 workshops over the next year
  - 1: Santa Fe, April 7-8
  - 2: France, week of July 20th
  - 3: Japan in the early Fall
- Broad engagement by the community
- Initial reports in summer 2009
- Final report for first year at SC09
- Planning for *IMMEDIATE* payoff
  - Could begin ramping up next year
- Stay tuned...



## Main Page

[Page](#) [Discussion](#) [View source](#) [History](#)

The mission of the **International Exascale Software Project (IESP)** is to lay the foundation for exascale computing by mobilizing the global open source software community to combine and coordinate their collective efforts far more efficiently and effectively than ever before. The IESP will hold a series of three workshops to organize and structure this community wide effort. The first, invitation-only workshop will occur on April 7th and 8th in Sante Fe, New Mexico, US, with people arriving in time for a reception on April 6th. Attendees will include members from industry, academia, and government, with expertise in a range of critical areas.

Goals for the first meeting include the following:

- ▶ Assess the short-term, medium-term and long-term needs of applications for peta/exascale systems
- ▶ Explore how laboratories, universities, and vendors can work together on coordinated HPC software
- ▶ Understand existing R&D plans addressing new programming models and tools addressing extreme scale, multicore, heterogeneity and performance
- ▶ Start development of a roadmap for software on extreme-scale systems

**Attendance at the workshop is by invitation only.** Additional details on registration will be coming soon.

### Workshop Information

- [Workshop Location](#)
- [Workshop Agenda \(draft\)](#)
- [Executive Committee](#)
- [Organizing Committee](#)
- [Background Material](#)