# *Overview of the Argonne Leadership Computing Facility*

*Scott Parker*

# *Argonne Leadership Computing Facility*

- ALCF was established in 2006 at Argonne to provide the computational science community with a leading-edge computing capability dedicated to breakthrough science and engineering

- One of two DOE national Leadership Computing Facilities (the other is the National Center for Computational Sciences at Oak Ridge National Laboratory)

- Supports the primary mission of DOE's Office of Science Advanced Scientific Computing Research (ASCR) program to discover, develop, and deploy the computational and networking tools that enable researchers in the scientific disciplines to analyze, model, simulate, and predict complex phenomena important to DOE.

# DOE INCITE Program
*Innovative and Novel Computational Impact on Theory and Experiment*

- **Solicits large computationally intensive research projects**
  - To enable high-impact scientific advances
  - Call for proposal opened once per year (call closed 6/30/2010)
  - INCITE Program web site: www.er.doe.gov/ascr/incite
- **Open to all scientific researchers and organizations**
  - Scientific Discipline Peer Review
  - Computational Readiness Review
- **Provides large computer time & data storage allocations**
  - To a small number of projects for 1-3 years
  - Academic, Federal Lab and Industry, with DOE or other support
- **Primary vehicle for selecting principal science projects for the Leadership Computing Facilities** (*60% of time at Leadership Facilities*)
- **In 2010, 35 INCITE projects allocated more than 600M CPU hours at the ALCF**

# DOE ALCC Program
## ASCR Leadership Computing Challenge

- Allocations for projects of special interest to DOE with an emphasis on high risk, high payoff simulations in areas of interest to the departments energy mission (30% of time at Leadership Facilities)
- Awards granted in June (call close on 2/15/2010)
  - http://www.er.doe.gov/ascr/facilities/alcc.htm
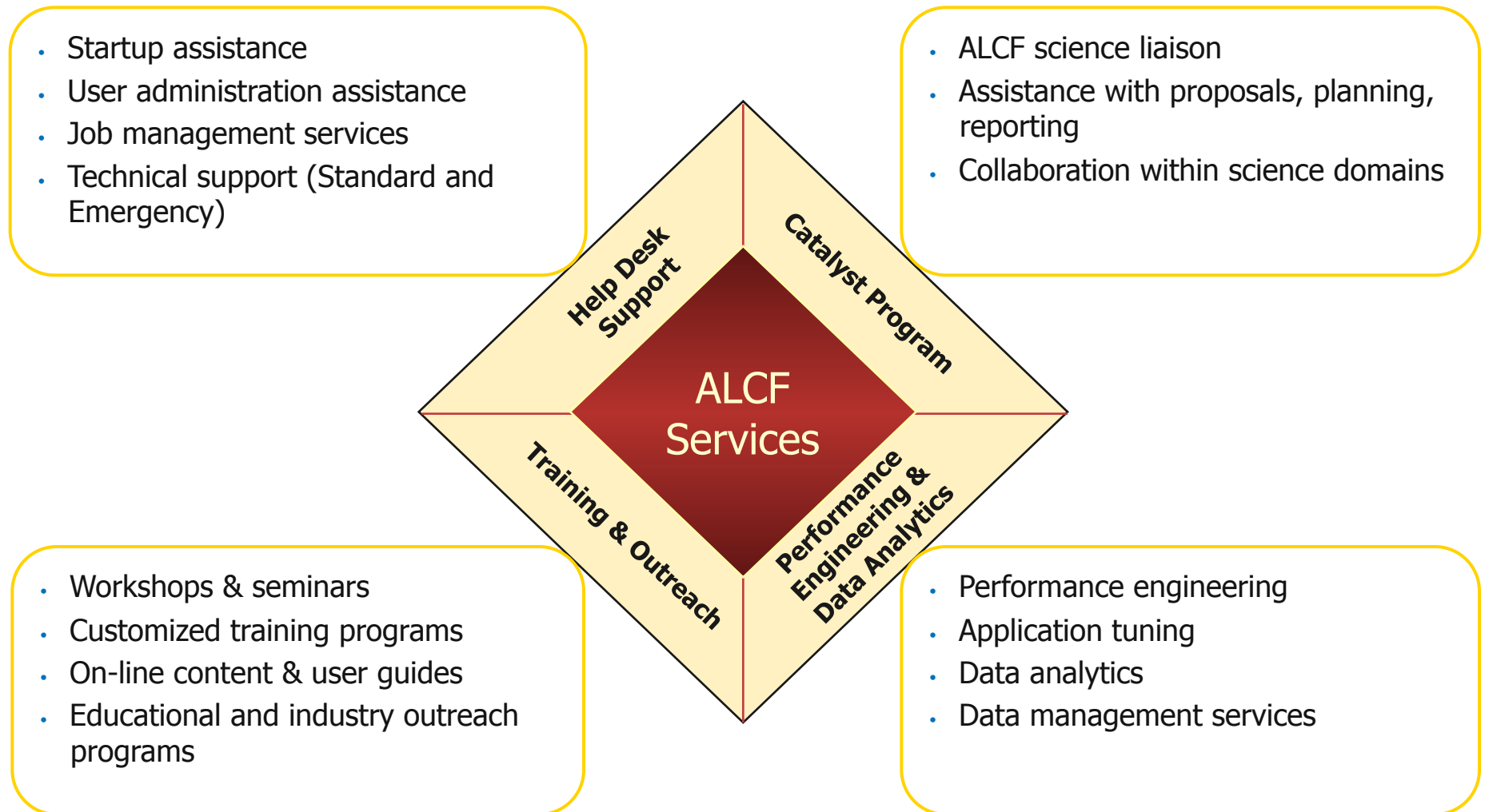- 9 awards at ALCF in 2010 for 300+ million core hours

# *Discretionary Allocations*

- Time is available for projects without INCITE or ALCC allocations!
- ALCF Discretionary allocations provide time for:
  - Porting, scaling, and tuning applications
  - Benchmarking codes and preparing INCITE proposals
  - Preliminary science runs prior to an INCITE award
- To apply go to the ALCF allocations page
  - www.alcf.anl.gov/support/gettingstarted

# *Get An Account!*

- If you don't have an account on an ALCF BG/P system (*Intrepid or Surveyor)* you can apply for a workshop account
- To apply:
  - Go to the URL: https://accounts.alcf.anl.gov/accounts/request.php
  - Select "Proceed with Account Request" at the bottom of the page
  - Select the project 'CScADS'
  - Foreign nationals require 593 forms which can take a while
- Running under the 'CScADS' project
  - User with account can run using the 'CScADS' project
    - *qsub –A CScADS  …*

# ALCF Service Offerings

- Startup assistance
- User administration assistance
- Job management services
- Technical support (Standard and Emergency)

- ALCF science liaison
- Assistance with proposals, planning, reporting
- Collaboration within science domains

**Help Desk Support**

**Catalyst Program**

## ALCF Services

**Training & Outreach**

**Performance Engineering & Data Analytics**

- Workshops & seminars
- Customized training programs
- On-line content & user guides
- Educational and industry outreach programs

- Performance engineering
- Application tuning
- Data analytics
- Data management services

# Argonne Leadership Computing Facility

- *Intrepid* - ALCF Blue Gene/P System:
  - 40,960 nodes / 163,840 PPC cores
  - 80 Terabytes of memory
  - Peak flop rate: 557 Teraflops
  - Linpack flop rate: 450.3
  - #9 on the Top500 list
- *Eureka* - ALCF Visualization System:
  - 100 nodes / 800 2.0 GHz Xeon cores
  - 3.2 Terabytes of memory
  - 200 NVIDIA FX5600 GPUs
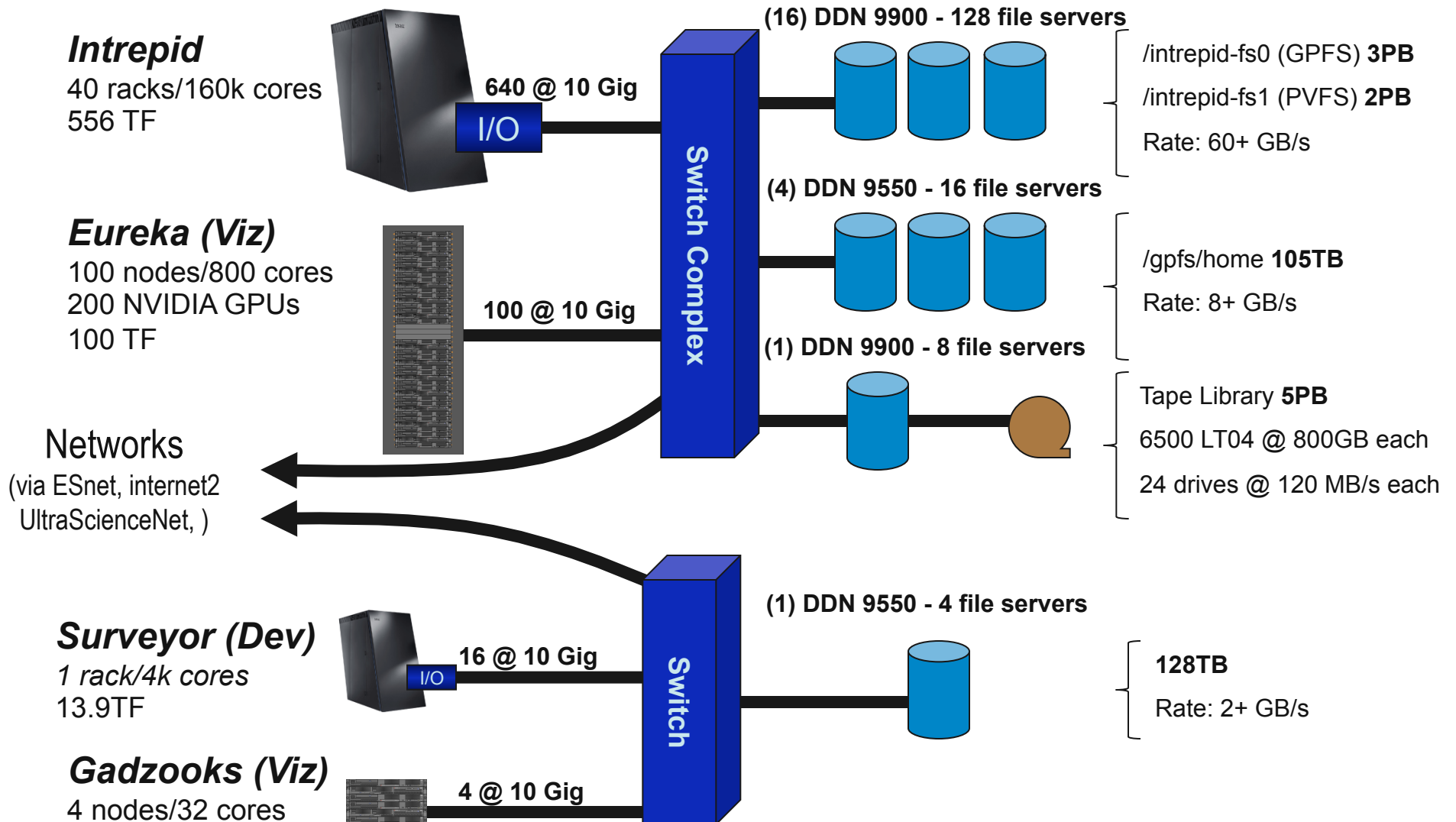  - Peak flop rate: 100 Teraflops
- Storage:
  - 6+ Petabytes of disk storage with an I/O rate of 80 GB/s
  - 5+ Petabytes of archival storage (10,000 volume tape archive)

# *The Next Generation ALCF System: BG/Q*

- Coming in early 2012 "Mira" a 10 Petaflop Blue Gene/Q system. An evolution of the Blue Gene architecture with:
  - Over 750k cores
  - 16 cores/node
  - 1 GB of memory per core
  - 70 PB of disk with 470 GB/s of I/O bandwidth
  - Power efficient, water cooled
- Argonne has worked closely with IBM over the last few years to help develop the specifications for this next generation Blue Gene system
- Early Science call concluded, allocations to be made soon
- Applications running on the BG/P should run immediately on the BG/Q, but may see better performance by exposing greater levels of parallelism at the node level

# ALCF Resources - Overview

**Intrepid**
40 racks/160k cores
556 TF

640 @ 10 Gig

**(16) DDN 9900 - 128 file servers**

/intrepid-fs0 (GPFS) **3PB**
/intrepid-fs1 (PVFS) **2PB**
Rate: 60+ GB/s

Switch Complex

**Eureka (Viz)**
100 nodes/800 cores
200 NVIDIA GPUs
100 TF

100 @ 10 Gig

**(4) DDN 9550 - 16 file servers**

/gpfs/home **105TB**
Rate: 8+ GB/s

**(1) DDN 9900 - 8 file servers**

Tape Library **5PB**
6500 LT04 @ 800GB each
24 drives @ 120 MB/s each

Networks
(via ESnet, internet2
UltraScienceNet, )

**(1) DDN 9550 - 4 file servers**

**Surveyor (Dev)**
*1 rack/4k cores*
13.9TF

16 @ 10 Gig

I/O

Switch

**128TB**
Rate: 2+ GB/s

**Gadzooks (Viz)**
4 nodes/32 cores

4 @ 10 Gig

# Blue Gene/P at ALCF

# BG/P: Covers removed

# BlueGene/P Overview

- 4 850Mhz PowerPC cores per chip
- 1 chip, 2 GB of DDR SDRAM, 5 network interfaces per compute node
- 32 compute nodes per node card
- 32 node cards per rack
- 1,024 nodes total per rack
- 40 rack on Intrepid

**Rack**
32 Node Cards
1024 chips, 4096 procs

**Intrepid System**
40 Racks

556 TF/s
82TB

14 TF/s
2 TB

**Node Card**
(32 chips  4x4x2)
32 compute, 0-2 IO cards

435 GF/s
64 GB

**Compute Card**
1 chip, 20 DRAMs

13.6 GF/s
2.0 GB DDR
Supports 4-way SMP

**Chip**
4 processors

850 MHz
8 MB EDRAM

**Front End Node / Service Node**
System p Servers
Linux SLES10

# *Blue Gene DNA*

- **Low power design** ⟹ massive parallelism
  - The leader in Green Computing
- **System on a Chip (SoC)**
  - Improves Price / Performance
  - Reduces system complexity & power
- **Custom designed ASIC**
  - Reducing overall part count, reducing errors
  - Permits tweaking CPU design to reduce soft errors
- **Dense packaging**
- **Fast communication network(s)**
- **Sophisticated RAS (reliability, availability, and serviceability)**
- **Dynamic software provisioning and configuration**
  - **Key feature for linking computer science with applications**

**3 watts per sustained gigaflops**

# Blue Gene/P ASIC

Argonne National Laboratory

# Double Hummer Floating Point Unit

8 Bytes      8 Bytes

Quad word Load
16 Bytes per instruction

Primary Registers    Secondary Registers

Full range of parallel and cross SIMD floating-point instructions

Quad word Store
16 Bytes per instruction

8 Bytes      8 Bytes

Quad word load/store operations require data aligned on 16-Byte boundaries.

# Blue Gene/P Interconnection Networks



- **3 Dimensional Torus**
  - Interconnects all compute nodes
  - Communications backbone for point-to-point
  - 3.4 Gb/s on all 12 node links (5.1 GB/s per node)
  - 0.5 µs latency between nearest neighbors, 5 µs to the farthest
  - MPI: 3 µs latency for one hop, 10 µs to the farthest
  - *Requires half-rack or larger partition*
- **Collective Network**
  - One-to-all broadcast functionality
  - Reduction operations for integers and doubles
  - 6.8 Gb/s of bandwidth per link per direction
  - Latency of one way tree traversal 1.3 µs, MPI 5 µs
  - Interconnects all compute nodes and I/O nodes
- **Low Latency Global Barrier and Interrupt**
  - Latency of one way to reach 72K nodes 0.65 µs, MPI 1.6 µs
- **10 Gb/s functional Ethernet**
  - Disk I/O
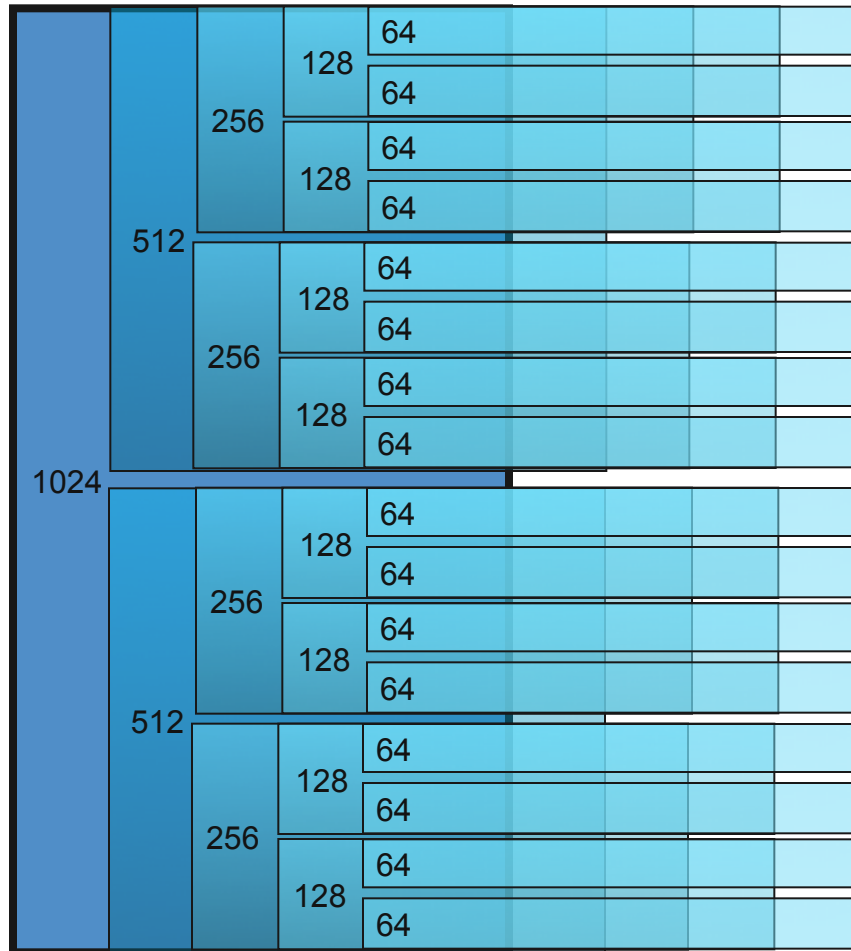- **1Gb private control (JTAG)**
  - Service node/system management

# Blue Gene/P Heterogeneity



- **Front-end nodes (FN)**: dedicated for user's to login, compile programs, submit jobs, query job status, debug applications, 2.5 GHz PowerPC 970, Linux OS
- **Service nodes (SN)**: perform system management services, create and monitoring processes, initialize and monitor hardware, configure partitions, control jobs, store statistics
- **I/O nodes (IO)**: provide a number of OS services, such as files, sockets, process management, debugging, 1 I/O node per 64 compute nodes
- **Compute nodes (CN)**: run user application, accessed only through qsub command, no shell, limited OS services, quad core 850 MHz PowerPC 450, CNK OS

# BG/P Partitions



Partitions on 1 rack of *Intrepid* showing number of nodes

- *Intrepid* compute nodes are grouped into partitions ranging from 64 to 40,960 nodes in powers of 2
- Jobs run in smallest partition into which they fit
- Job makes entire partition unavailable
- Only 1 job may run in a partition
- Smaller partitions are enclosed inside of larger ones
- Minimum partition size is 64 nodes
  - *1 I/O node for each 64 compute nodes*
- Partition's networks electrically isolated, each partition is it's own torus/mesh
- Partitions <512 nodes form a mesh network, partitions >=512 nodes form a torus

# *Blue Gene/P Software*

- System:
  - Linux on Login and I/O Nodes
  - Compute Node Kernel (CNK) O/S on Compute Nodes (Linux like)
  - XL compilers (C, C++, Fortran 77-90-95-2003)
  - Python
  - MPI/OpenMP
  - ESSL math libraries
- Management:
  - BlueGene/P Control system (IBM DB2 database)
  - Cobalt – resource manager(qsub, qstat, qdel, qalter)
  - Clusterbank – allocation management system
- Storage:
  - GPFS – parallel filesystem
  - PVFS – high performance parallel filesystem
  - Tape systems - HPSS, Amanda

# *Unique and Challenging Features*

- Low power cores (850 MHz, 3.4 GFlop) but many of them (163,840)
  - High scalability is key to high total flop rate
  - Balance between CPU and network makes scalability easier
- Relatively low memory per CPU core (but very large aggregate).
  - BG/P: 2 Gig / 4 cores
  - True SMP is possible (sharing data structures)
- Single node optimization
  - Use of "double hummer" requires lots of hand tuning and compiler experimentation
  - Strategies: Good math libs, performance counters, code tools
- Scalable I/O strategy required
  - One file per process **strongly** discouraged
  - PnetCDF and HDF5 are good strategies for effectively using parallel storage system (GPFS, PVFS)
- Debugging at scale remains challenging
  - Tool groups are helping, but the issue is nevertheless hard

# *More about BlueGene/P*

- Logical partitions, with complete electrical isolation
- Partition rebooted between jobs
- Only allowed limited thread/process per core using one of three modes (1 thread/process per core is generally optimal)
    - SMP – 1 process with 4 threads - 1 thread per core
    - Dual – 2 processes with 2 thread
    - VN – 4 processes – 1 process per core
- Specialized IBM kernels – compute node kernel, I/O node kernel
    - Single-executable kernel for compute nodes
    - Usually runs MPI code, can also run in "HTC" mode
    - Stripped-down Linux for I/O nodes
    - The ciod daemon handles system calls
- But also can run custom kernels
    - *ZeptoOS / Compute Node Linux*
    - *ZeptoOS I/O node kernel*
    - *Plan 9 (INCITE project)*

# *Programming Models and Development Environment*

- Languages:
  - Full language support with IBM XL and GNU compilers
  - Languages: Fortran, C, C++, Python
- MPI:
  - Based on MPICH2 1.0.x base code:
    - *MPI-IO supported*
    - *One-sided communication supported*
    - *No process management (MPI_Spawn(), MPI_Connect(), etc)*
  - Utilizes the 3 different BG/P networks for different MPI functions
- Threads:
  - OpenMP 2.5
  - NPTL Pthreads
- Linux development environment:
  - Compute Node Kernel provides look and feel of a Linux environment
    - *POSIX routines (with some restrictions: no fork() or system())*
    - *BG/P adds pthread support, additional socket support*
  - Supports statically and dynamically linked libraries (static is default)
  - Cross compile since login nodes and compute nodes have different processor & OS

# *Restrictions and Complications*

- SPMD model:
  - By default compute nodes run the same executable (a few work-arounds available)
- Space Sharing:
  - By default one parallel job per partition of machine
  - Optimized for one process/thread per core in each compute node
    - *smp-mode, one MPI task/node, 4 threads/task, 2GB of RAM*
    - *dual-mode, two MPI tasks/node, 2 threads/task, 1GB of RAM*
    - *vn-mode, 4 single-threaded MPI tasks/node, 512MB of RAM*
- memory limited to physical memory (no virtual memory)
- restricted set of POSIX routines (no fork, system, …)
- Cross compiling required due to hardware & O/S differences between login and compute nodes

# *Performance and Debugging Tools*

- IBM High Performance Computing Toolkit
  - MPI Profile and Tracing Library
  - HPM Library for hardware performance counters
  - Xprofiler visualization of gprof profiles
- Rice HPCToolkit:
  - Sample based profiling of applications
- TAU – Tuning and Analysis Utilities
  - MPI Profiling
  - Performance counter dat
- Hardware Counters (Universal Performance Counters – UPC)
  - Blue Gene/P provides 256 on chip counters for hardware events
  - UPC and HPM libraries provide access to counters
- gprof – standard linux profiling tool
- Core files – lightweight core files, text format, no full memory dump
- Coreprocessor - Generates stack trace from core files
- GDB – gnu debugger
- TotalView - Debugging across ten of thousands of cores

# Supported Libraries and Programs

| Library | Location | Description |
| --- | --- | --- |
| BLAS, LAPACK | /soft/apps/blas-lapack-lib | Basic vector linear algebra subroutines. |
| ESSL | /soft/apps/ESSL-4.3 | Mathematical subroutines designed to improve the performance of engineering and scientific applications on BlueGene |
| LIBGOTO | /soft/apps/LIBGOTO | Very efficient BLAS-1.2.3 implementation for BlueGene from Kazushige |
| SCALAPACK | /soft/apps/SCALAPACK | High-performance linear algebra routines for distributed-memory message-passing MIMD computers and networks of workstations. |
| PETSc | /soft/apps/petsc | A suite of data structures and routines that provide the building blocks for the implementation of large-scale application codes on serial and parallel computers. |
| fftw-2.1.5, fftw-3.1.2 | /soft/apps/fftw- | A library for computing the discrete many-dimensional Fourier transform |
| p3dfft | /soft/apps/p3dfft-2.1beta- | Highly scalable parallel 3D Fast Fourier Transforms library. |
| hypre-2.0.0 | /soft/apps/hypre-2.0.0 | A library for solving large, sparse linear systems of equations on massively parallel computers. |
| SuperLU-3.0 | /soft/apps/SuperLU_3.0 | A library for the direct solution of large, sparse, non-symmetric systems of linear equations on high performance machines. |
| MUMPs-4.7.3 | /soft/apps/MUMPS_4.7.3 | A multifrontal massively parallel sparse direct solver. |
| spooles-2.2 | /soft/apps/spooles-2.2 | A library for solving sparse real and complex linear systems of equations. |

# Supported Libraries and Programs

| Program | Location | Description |
|---------|----------|-------------|
| TotalView | /soft/apps/totalview-8.5.0-0 | Multithreaded, multiprocess source code debugger for high performance computing. |
| Coreprocessor | /soft/apps/coreprocessor.pl | A tool to debug and provide postmortem analysis of dead applications. |
| TAU-2.17 | /soft/apps/tau | A portable profiling and tracing toolkit for performance analysis of parallel programs written in Fortran, C++, and C |
| HPCT | /soft/apps/hpct_bgp | MPI profiling and tracing library, which collects profiling and tracing data for MPI programs. |

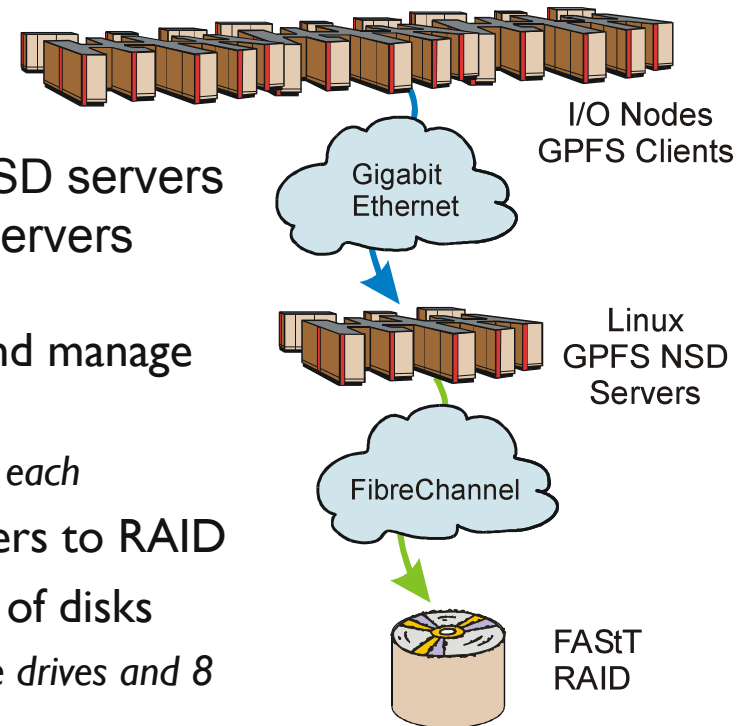| Program | Location | Description |
|---------|----------|-------------|
| armci | /bgsys/drivers/ppcfloor/comm | The Aggregate Remote Memory Copy (ARMCI) library |
| HDF5 | /soft/apps/hdf5-1.6.6 | The Hierarchical Data Format (HDF) is a model for managing and storing data. |
| NetCDF | /soft/apps/netcdf-3.6.2 | A set of software libraries and machine-independent data formats that supports the creation, access, and sharing of array-oriented scientific data. |
| Parallel NetCDF | /soft/apps/parallel-netcdf-1.0.2 | A library providing high-performance I/O while still maintaining file-format compatibility with Unidata's NetCDF. |
| mercurial-0.9.5 | /soft/apps/mercurial-0.9.5 | A distributed version-control system |
| Scons | /soft/apps/scons-0.97 | A cross-platform substitute for the classic Make utility |
| tcl-8.4.14 | /soft/apps/tcl-8.4.14 | A dynamic programming language, |

# File Systems

**Intrepid**
40 racks/160k cores
556 TF

**(16) DDN 9900 - 128 file servers**

**640 @ 10 Gig**

I/O

Switch Complex

/intrepid-fs0 (GPFS) **3PB**

/intrepid-fs1 (PVFS) **2PB**

Rate: 60+ GB/s

**(4) DDN 9550 - 16 file servers**

**Eureka (Viz)**
100 nodes/800 cores
50 NVIDIA S4 GPUs
100 TF

**100 @ 10 Gig**

/gpfs/home **100TB**

Rate: 8+ GB/s

**(1) DDN 9900 - 8 file servers**

Networks
(via ESnet, internet2
UltraScienceNet, )

Tape Library **5PB**

6500 LT04 @ 800GB each

24 drives @ 120 MB/s each

**(1) DDN 9550 - 4 file servers**

**Surveyor (Dev)**
*1 rack/4k cores*
13.9TF

I/O

**16 @ 10 Gig**

Switch

**128TB**

Rate: 2+ GB/s

**Gadzooks (Viz)**
4 nodes/32 cores

**4 @ 10 Gig**

# Intrepid File Systems

- /gpfs/home
  - GPS, 105 TB
  - Intended for storing source, configuration, and input files
  - Not intended for larger scale parallel I/O, or storing large files
  - Backup and snap-shot files
- /intrepid-fs0
  - GPFS, 3 PB, 60+ GB/s
  - Intended for very fast parallel IO, program input and output
  - Not backed up, but you can initiate archive via HPSS
- /intrepid-fs1
  - PVFS, 2.2 PB, 50+ GB/s
  - Intended for very fast parallel IO, program input and output
  - Not backed up, but you can initiate archive via HPSS
  - NOTE: Binaries can not be executed from PVFS

# General Parallel File System (GPFS) on Blue Gene

- Blue Gene can generate enormous I/O demand
    - BG/P IO-rich has 640 I/O nodes at 10Gb/s
    - Requires a parallel file system (peak 78 GB/s)
- GPFS Setup:
    - **I/O nodes** run GPFS client that call external NSD servers
    - **10 GB/s Ethernet** connect I/O nodes to NSD servers
        - *900+ port 10 Gigabit Ethernet Myricom switch complex*
    - **NSD Servers** run parallel file system software and manage incoming FS traffic from I/O nodes
        - *136 two dual core Opteron servers with 8 Gbytes of RAM each*
    - **Infiniband/Fiber Channel** connect NSD servers to RAID
    - **Enterprise storage** controllers and large racks of disks
        - *17 DataDirect S2A9900 controller pairs with 480 1 Tbyte drives and 8 InfiniBand ports per pair*
- Brings traditional benefits of GPFS to Blue Gene
    - I/O parallelism
    - Cache consistent shared access
    - Aggressive read-ahead, write-behind

I/O Nodes
GPFS Clients

Gigabit
Ethernet

Linux
GPFS NSD
Servers

FibreChannel

FAStT
RAID

**Argonne National Laboratory**

# *Visualization and Data Analytics*

- *Eureka* makes data analytics and visualization at Intrepid's scale possible through the world's largest installation of NVIDIA S4 external GPUs

- The system consists of 100 servers with 200 Quadro FX5600 graphics engines.

- The system is attached directly to the core switch complex of the Blue Gene/P, providing very high throughput between the BG/P and the analysis nodes, and also to the parallel file system

# *Eureka Visualization System*

- **100 Nodes:**
  - Two 2.0 GHz Quad Core Xeon (8 cores/node)
  - 32 GB RAM
- **50 NVidia S4 external GPUs**
  - 200 Quadro FX5600 high end graphics cards
- **Nodes connect to the S4 via a 16x PCIe V2.0 card**
- **Over 111 TF peak FLOP rating**
  - *Includes the GPUs*
- **3.2 TB of RAM (5% of intrepid RAM)**
- **No local scratch space**
  - Data access is all via the central parallel file system
  - Myricom 10G (10 Gbs) NIC

# Software Resources

■ Installed and Supported Software:

– VisIT

– ParaView

– VTK

– VMD

– And good old gnuplot

■ Software stack is driven by our users

– So if we don't have what you need, please speak up

– Note that we don't list any commercial apps

• *No user driven demand to date*

# *ALCF INCITE Projects Span Many Domains*

Life Sciences
U CA-San Diego

Applied Math
Argonne Nat'l Lab

Physical Chemistry
U CA-Davis

Nanoscience
Northwestern U

Engineering Physics
Pratt & Whitney

Biology
U Washington

# 2009 INCITE Allocations at ALCF

**35projects**

**600+ M CPU Hours**

# *Final Thoughts on Blue Gene*

- General purpose architecture that excels in virtually all areas of computational science
- Though having many special features it presents an essentially standard Linux/PowerPC programming environment
- Pursues performance through high levels of parallelism with good balance between processor, memory, and network speed
- Significant impact on HPC – 3 of the top 10 machines are currently Blue Gene systems
- Next Generation Blue Gene system is on its way – Lawrence Livermore and Argonne have plans for the acquisition of 10-20 Petaflop Blue Gene/Q systems
- Delivers excellent performance per watt, performance per square foot
- High reliability and availability
- Able to run applications consistently and with high performance across the entire system

# *If you want to know more…*

- ALCF web site: www.alcf.anl.gov
  - Information on ALCF system and activities
  - Information on applying for accounts
- Getting Started Guide:
  - https://wiki.alcf.anl.gov/index.php/Quick_Reference_Guide
- ALCF Support Wiki: wiki.alcf.anl.gov
  - Documentation
  - FAQ
- Support email address: support@alcf.anl.gov
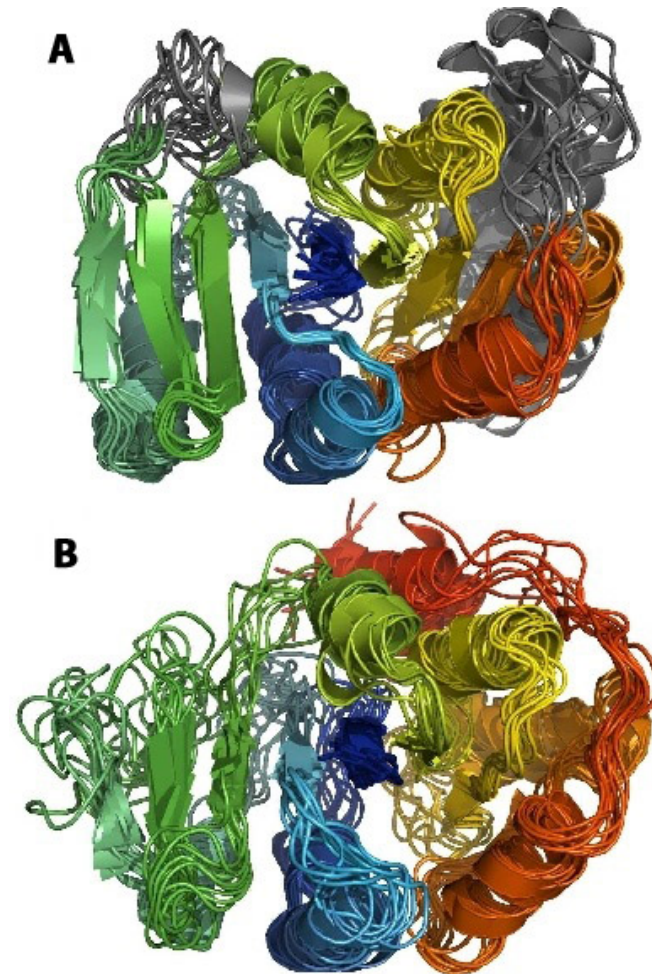  - Any question: account, technical, etc.

# *Science Projects Using ALCF Resources*

- Computational Protein Structure Prediction and Design
- Heat/fluid dynamics for advanced burner reactor design
- Gating Mechanism of Membrane Proteins
- Lattice Quantum Chromodynamics
- Validation of Type Ia supernova models
- Next-generation Community Climate System Model
- Simulation of Two Phase Flow Combustion in Gas Turbines
- Probing the Non-Scalable Nano Regime in Catalytic Nanoparticles
- Reactive MD Simulation of Nickel Fracture

# Computational Protein Structure Prediction and Protein Design

*David Baker, University of Washington*

- Code: Rosetta
- Rapidly determine an accurate, high-resolution structure of any protein sequence up to 150-200 residues
- Incorporating sparse experimental NMR data into Rosetta to allow larger proteins



*Comparison of computational and experimentally derived protein structure*

nking of all CASP8 groups

# Computational Nuclear Engineering

*Paul Fischer, Argonne National Lab*

- Code: Nek5000

- Improve the Safety, Efficiency and Cost of Liquid-Metal-Cooled Fast Reactors through computation

- High resolution 3D simulation of fluid flow and heat transfer in full reactor core sub-assembly

- Simulation requires full pin assembly, cannot use symmetries
  - 217 Pin configuration
  - 2.95 million spectral elements
  - 1 billion grid points

# Gating Mechanism of Membrane Proteins

*Benoit Roux, ANL, University of Chicago*

- Code: NAMD with periodicity and particle-mesh Ewald method
- Understand how proteins work so we can alter them to change their function
- Validated the atomic models of Kv1.2 and first to calculate the gating charge in the two functional states
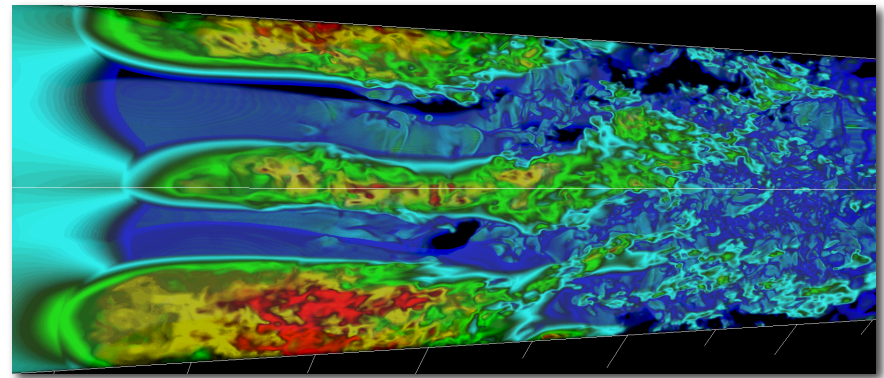
# Lattice QCD

*Bob Sugar and US-QCD*

- Code: MILC
- Determine parameters for Standard Model, including quark mass
- Addresses fundamental questions in high energy and nuclear physics
- Directly related to major experimental programs
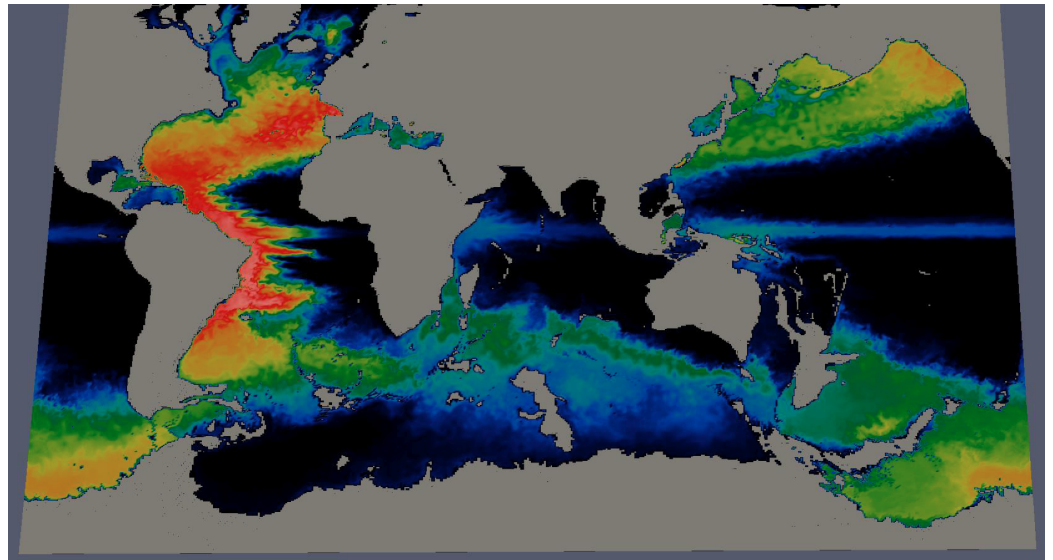
# FLASH

*Don Lamb, University of Chicago*

- Code: FLASH
- Investigating Type 1a supernovae which are used as a "standard candle" is astronomy
- Answered critical question on critical process in Type Ia supernovae
  - First simulated buoyancy-driven turbulent nuclear combustion in the fully-developed turbulent regime while also simultaneously resolving the Gibson scale
  - Reveals complex structure
  - Requires higher resolution studies
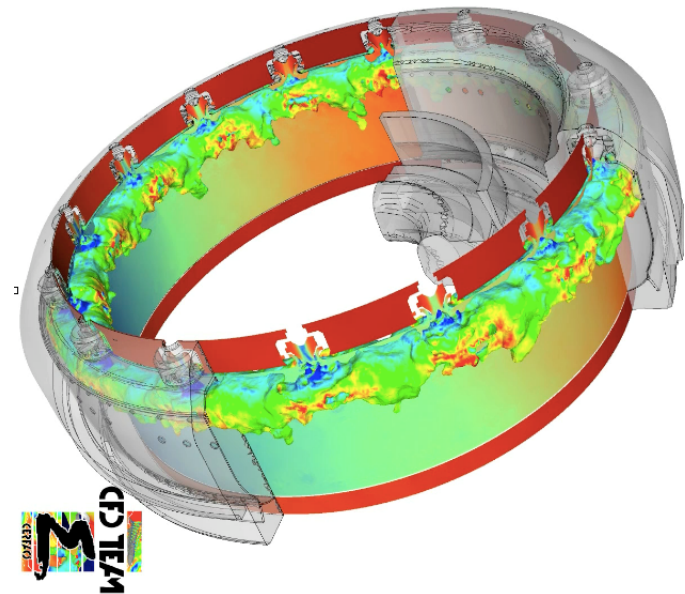
# Climate Simulations

## Warren Washington, NCAR

- Code: CCSM (climate simulation code used by the DOE and NSF climate change experiments)
- Ultra-high resolution atmosphere simulations:
  - CAM-HOMME atmospheric component
  - 1/8th degree (12.km avg grid) coupled with land model at 1/4th degree and ocean/ice
  - Testing up to 56K cores, 0.5 simulated years per day with full I/O
  - ½ degree finite-volume CAM with tropospheric chemistry and 399 tracer

# *Large Eddy Simulation of Two Phase Flow Combustion in Gas Turbines*
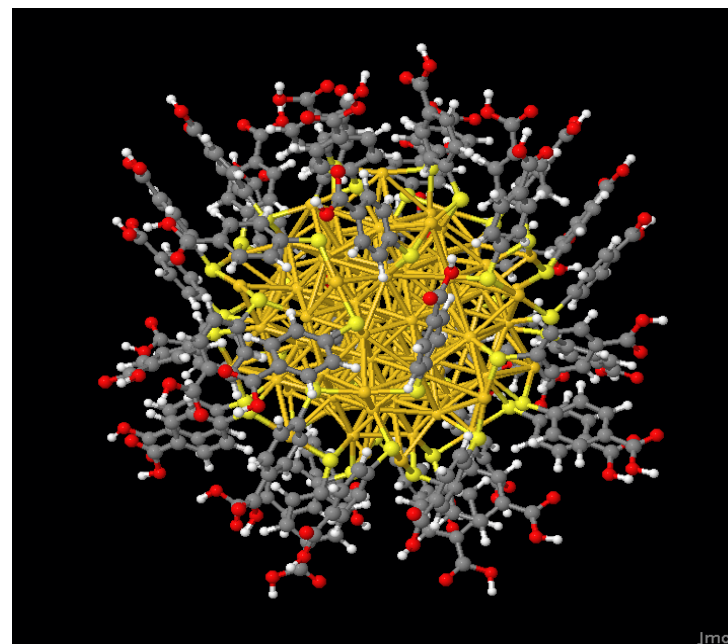
*Thierry Poinsot, CERFACS*

- Code: AVBP
- Improve airplane and helicopter engine designs – reduce design cost, improve efficiency
- Study the two phase flows and combustion in gas turbine engines
- Results:
  - First and only simulations of a full annular chamber
  - Identified instability mechanisms reducing efficiency in gas turbines

# Probing the Non-Scalable Nano Regime in Catalytic Nanoparticles

*Jeff Greeley, Argonne National Laboratory*

- Code: GPAW
- Gold has exciting catalytic capabilities that are still poorly understood
  - e.g. CO into $CO_2$
- First principle calculations of the actual nano-particles are needed

# Reactive MD Simulation of Nickel Fracture

*Priya Vashista, USC*

- Sulfur atoms in nickel grain boundaries embrittle the nickel
    - Phys. Rev. Letters 16 Apr 2010
- 65K cores on Intrepid BlueGene/P
    - 48 million atoms
    - Chemically reactive molecular dynamics
    - Fracture mechanics modeled



Embrittlement of Nickel: Sulfur segregation on grain boundaries

Pure Ni | With S

7.5nm | 7.5nm