# ORNL Cray XT3/4 Overview

**Jeffrey S. Vetter**

**And a cast of dozens …**

**Cray, NCCS, Pat Worley, Sadaf Alam, Weikuan Yu**

# Software Development Tools for Petascale Computing Workshop on Aug 1-2

# Current NCCS Resources

February 2007 Summary

7 Systems

Network Routers

Control Network
1 GigE
10 GigE
UltraScience

| CRAY XT3 JAGUAR | CRAY X1E PHOENIX | SGI ALTIX RAM | IBM SP4 CHEETAH | IBM LINUX NSTG | VISUALIZATION CLUSTER | IBM HPSS |
|---|---|---|---|---|---|---|
| (23016) 2.6GHz 46TB Memory | (1,024) 0.5GHz 2 TB Memory | (256) 1.5GHz 2TB Memory | (864) 1.3GHz 1.1TB Memory | (56) 3GHz 76GB Memory | (128) 2.2GHz 128GB Memory | Many Storage Devices Supported |
| 900 TB | 44 TB | 36 TB | 32 TB | 4.5 TB | 9 TB | 5 TB |

Supercomputers
25,344 CPUs
51 TB Memory
144 TFlops

Total Shared Disk
1.03 PB

10 PB

**Scientific Visualization Lab**

27 projector, 35 megapixel, Power Wall

**Test Systems**
• 96 processor Cray XT3
• 32 processor Cray X1E*
• 16 Processor SGI Altix

**Evaluation Platforms**
• 144 processor Cray XD1 with FPGAs
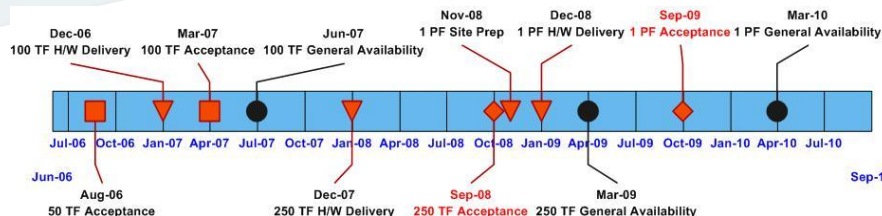• SRC Mapstation
• Clearspeed

**Backup Storage**

10PB

3 CSM Experimental Computing Lab houses IBM CELL, GPUs, Clearspeed, FPGAs, etc

# NCCS Roadmap for Leadership Computing

| Mission: Deploy and operate the computational resources needed to tackle global challenges | Vision: Maximize scientific productivity and progress on the largest scale computational problems |
|---|---|
| • **Future Energy**<br>• **Understanding the universe**<br>• **Nanoscale materials**<br>• **Climate Change**<br>• **Computational Biology** | • **Providing world class computational resources and specialized services**<br>• **Providing a stable hardware/software path of increasing scale to maximize productive applications development**<br>• **Work with users to scale applications to take advantage of systems** |

Cray XT4: 119 TF
Cray X1E: 18 TF

Cray XT4: 250 TF
Cray X1E: 18 TF

Cray HPCS-0: 1 PF leadership class system for science

Future leadership class sustained PF system for science

**FY2007**  **FY2008**  **FY2009**  **FY2011**

Dec-06 100 TF H/W Delivery
Mar-07 100 TF Acceptance
Jun-07 100 TF General Availability
Nov-08 1 PF Site Prep
Dec-08 1 PF H/W Delivery
Sep-09 1 PF Acceptance
Mar-10 1 PF General Availability

Jun-06
Aug-06 50 TF Acceptance
Dec-07 250 TF H/W Delivery
Sep-08 250 TF Acceptance
Mar-09 250 TF General Availability
Sep-10

Jul-06 Oct-06 Jan-07 Apr-07 Jul-07 Oct-07 Jan-08 Apr-08 Jul-08 Oct-08 Jan-09 Apr-09 Jul-09 Oct-09 Jan-10 Apr-10 Jul-10

# Jaguar – Cray XT4 with 11,706 Dual-Core AMD Opteron Processors

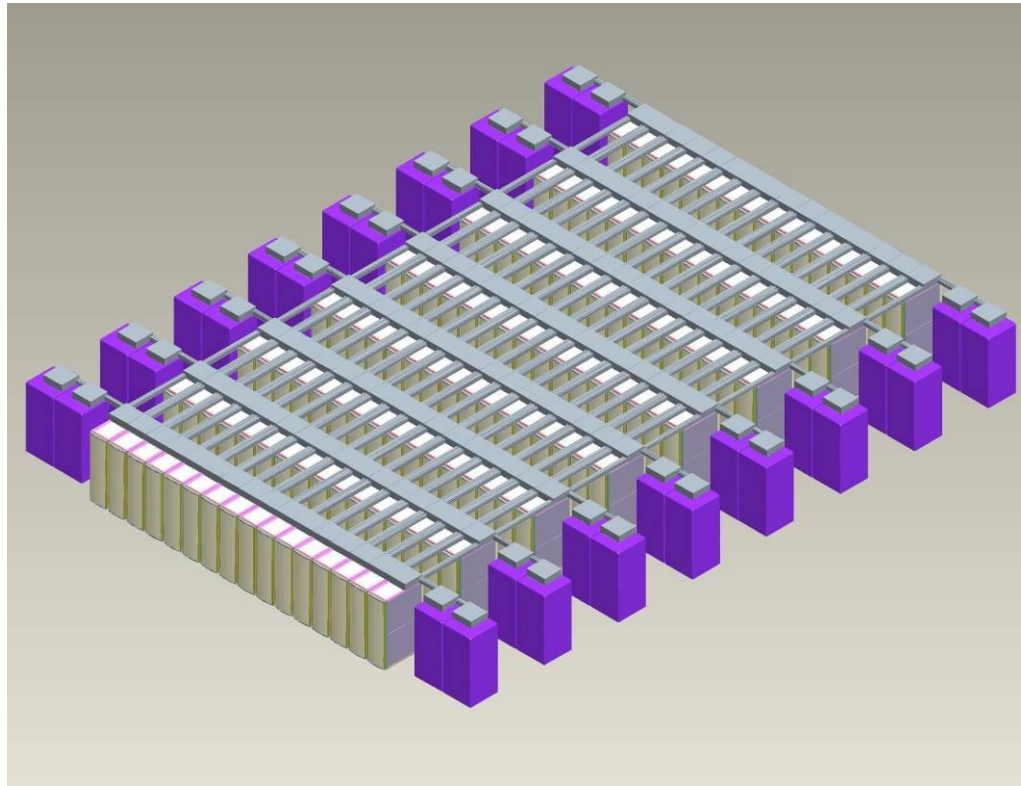- 119 TF peak performance
- 46 TB main memory
- 2.5 MW of power (20 KW per rack)
- 124 cabinets – 96 Opterons per cabinet
- Air cooled bottom-to-top
- Single 3000 CFM variable speed fan per rack
- 2.3 miles of interconnect cables
- Upgrade to Quad-core processors in Fall, 2007

# 1000 TF Cray "Baker" system in 2008

System configuration
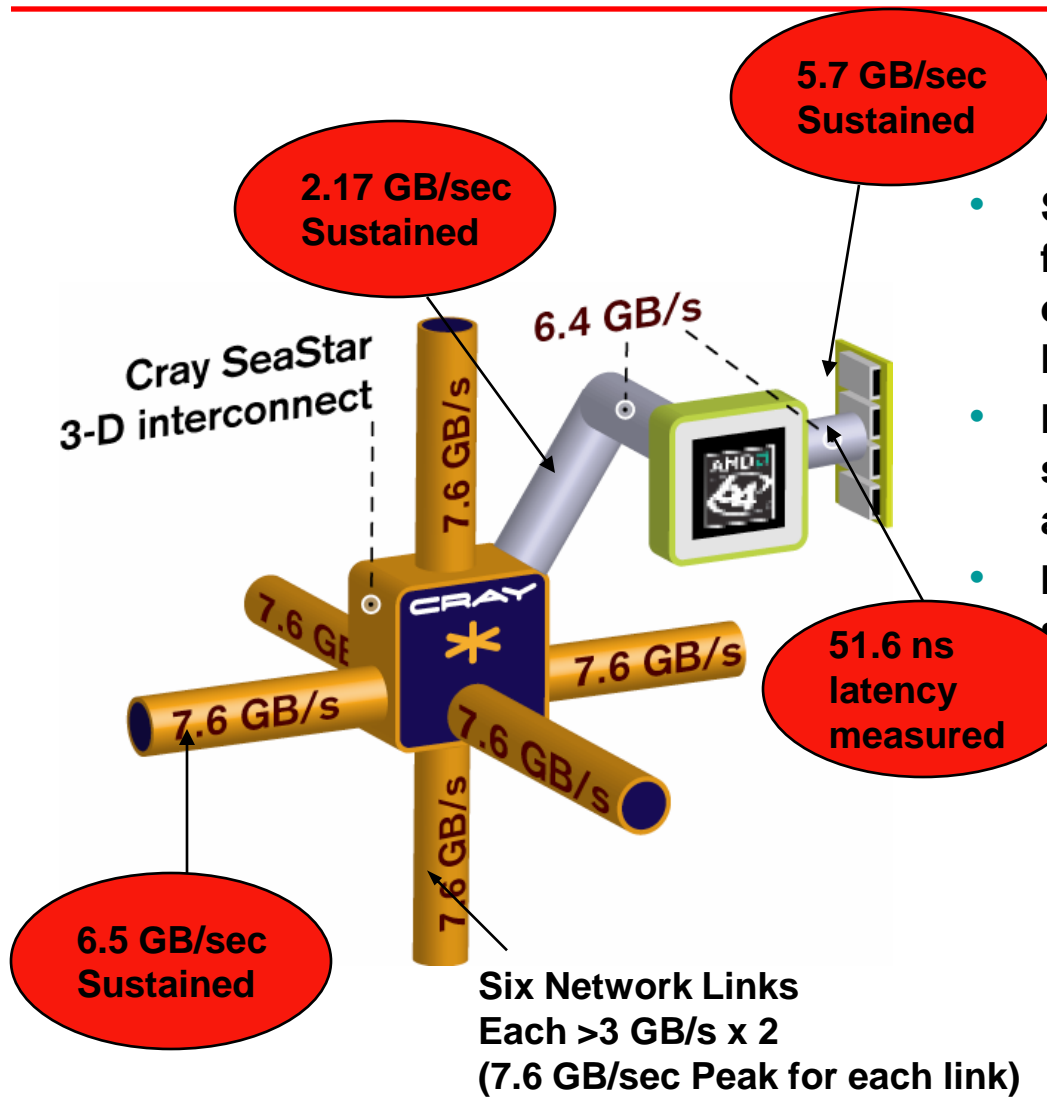
➡ 1 PF peak

➡ ~24,000 quad-core processors

➡ ~50 KW per cabinet

➡ ~7 MW power

➡ Over 3,000 watts/ft$^2$

➡ 40+ heat exchange units (chilled water to R-134a)



*1 PF Cray system in 2008*
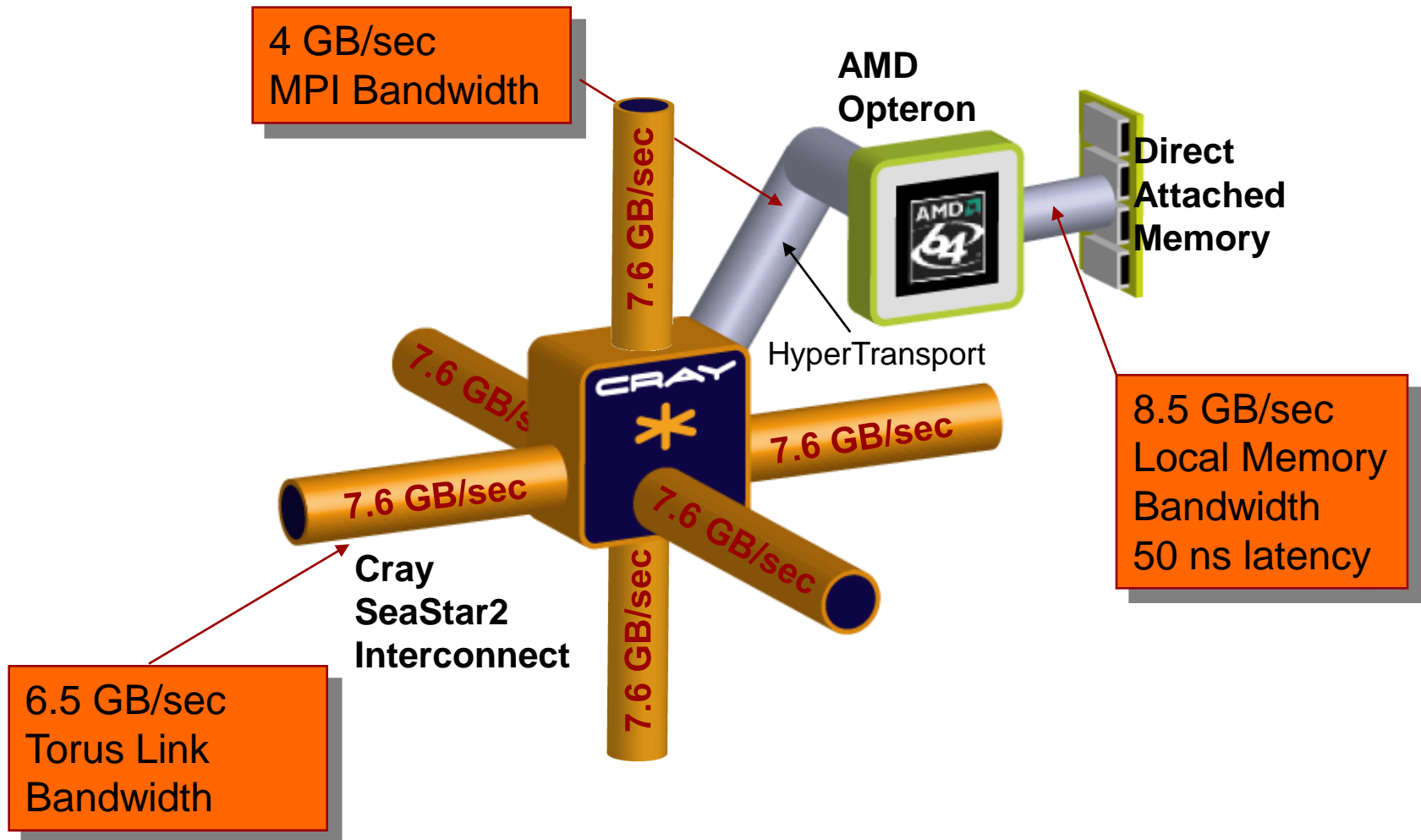
Used by permission: Cray, Inc.

# Cray XT3 Processing Element: Measured Performance

**5.7 GB/sec Sustained**

**2.17 GB/sec Sustained**

6.4 GB/s

Cray SeaStar 3-D interconnect

7.6 GB/s

7.6 GB/s

7.6 GB/s

7.6 GB/s

7.6 GB/s

7.6 GB/s

CRAY

**51.6 ns latency measured**

**6.5 GB/sec Sustained**

**Six Network Links
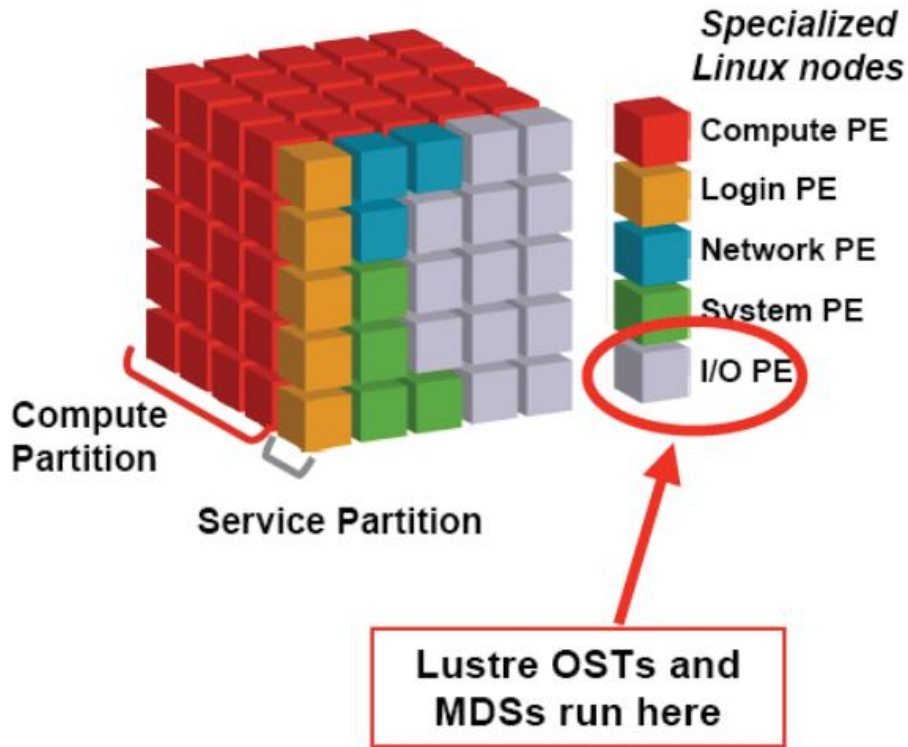Each >3 GB/s x 2
(7.6 GB/sec Peak for each link)**

- **SDRAM memory controller and function of Northbridge is pulled onto the Opteron die. Memory latency reduced to <60 ns**

- **No Northbridge chip results in savings in heat, power, complexity and an increase in performance**

- **Interface off the chip is an open standard (HyperTransport)**

# The Cray XT4 Processing Element: Providing a bandwidth-rich environment



4 GB/sec
MPI Bandwidth

AMD
Opteron

7.6 GB/sec

Direct
Attached
Memory

HyperTransport

7.6 GB/sec

7.6 GB/sec

7.6 GB/sec

7.6 GB/sec

7.6 GB/sec

8.5 GB/sec
Local Memory
Bandwidth
50 ns latency

Cray
SeaStar2
Interconnect

6.5 GB/sec
Torus Link
Bandwidth

# I/O

# I/O Configuration



**Specialized Linux nodes**

Compute PE
Login PE
Network PE
System PE
I/O PE

Compute Partition

Service Partition

Lustre OSTs and MDSs run here
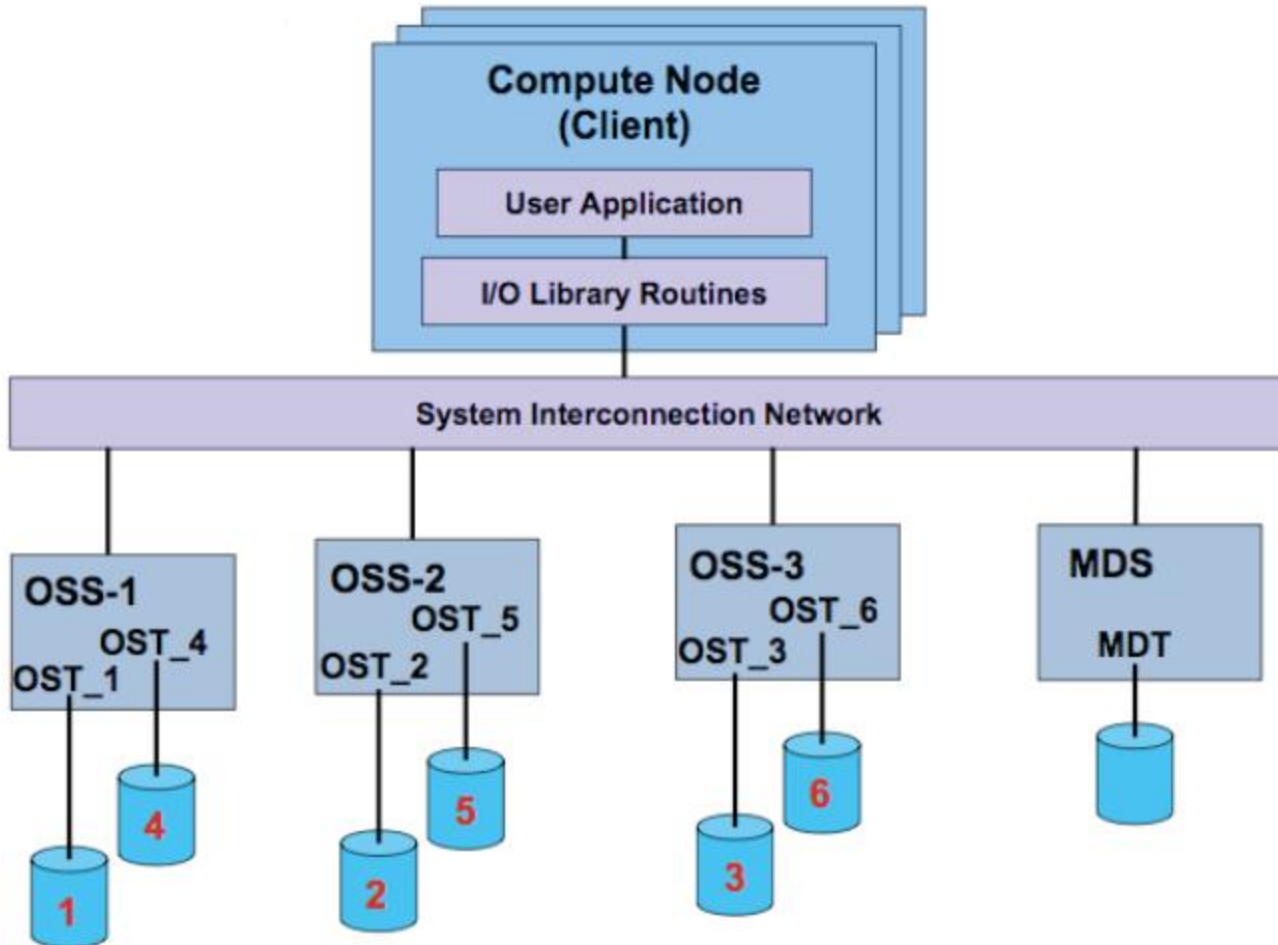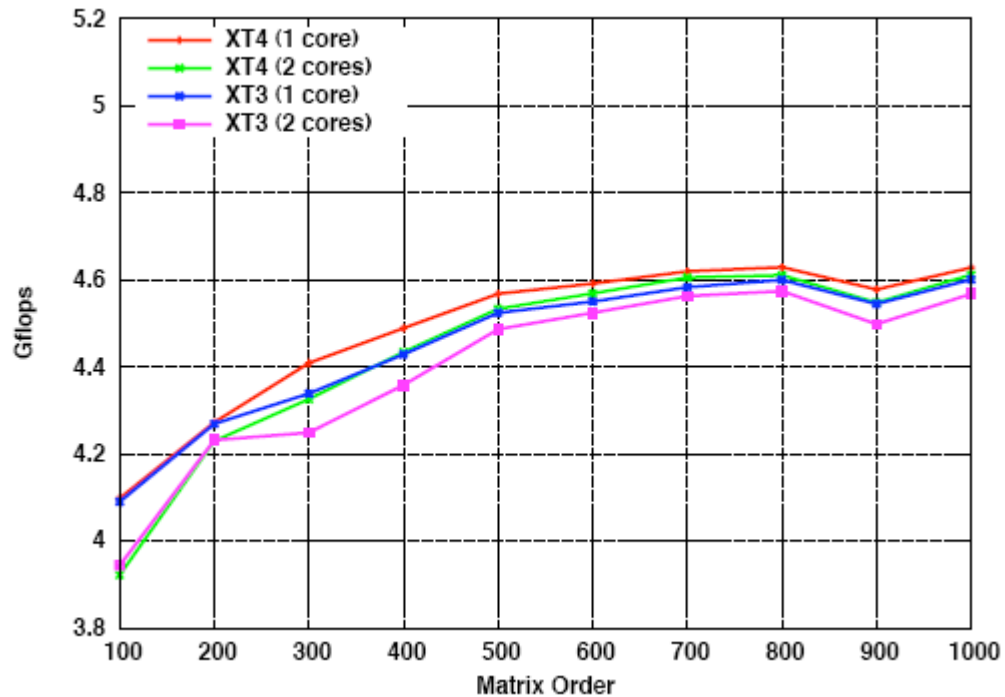
➡ Lustre filesystems

- Serviced by 80 I/O nodes
- /lustre/scr144
  - 144 OSTs
  - Peak 72 GB/s
  - Target ~48 GB/s
  - Early results
    - Read 45 GB/s
    - Write 25 GB/s
- /lustre/scr72[a,b]
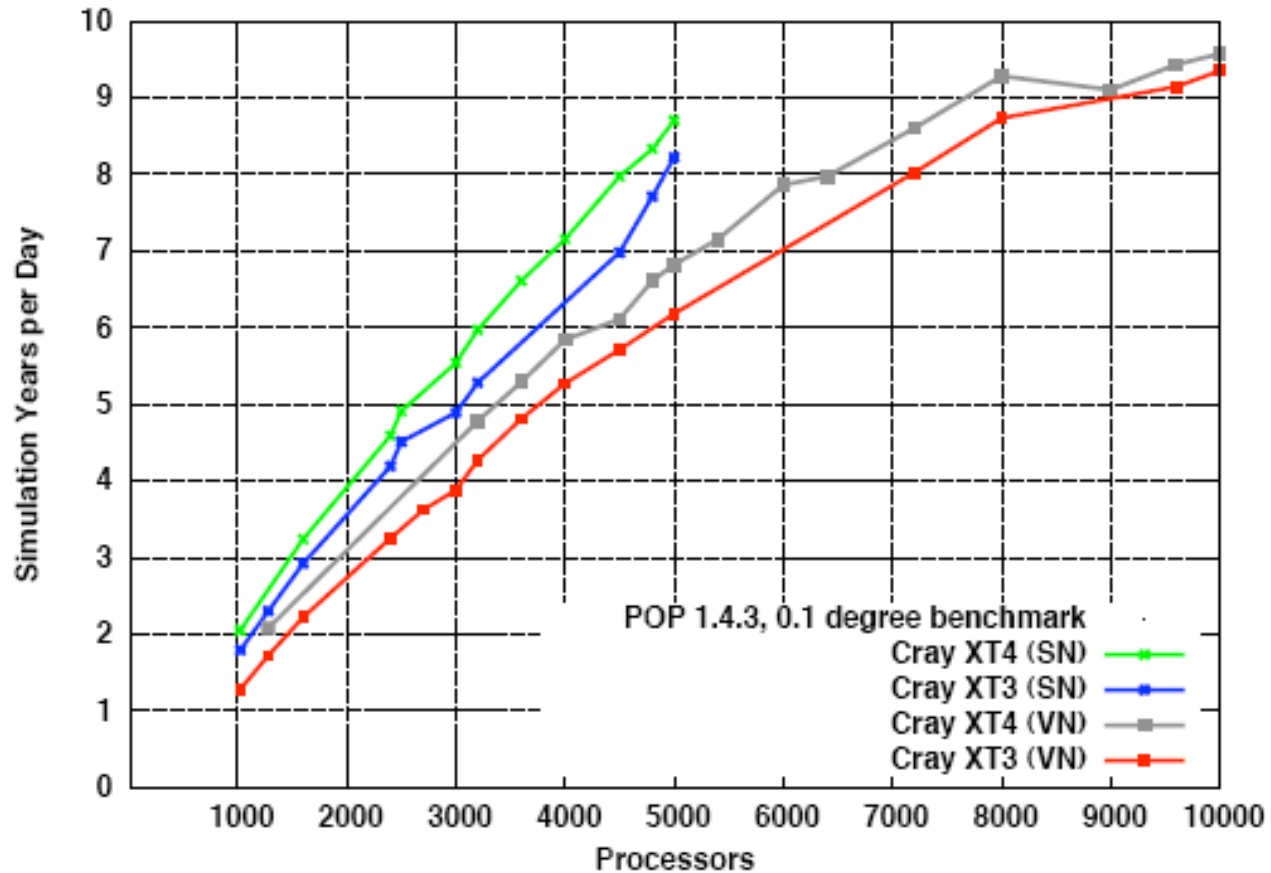  - 72 OSTs each
  - Default scratch
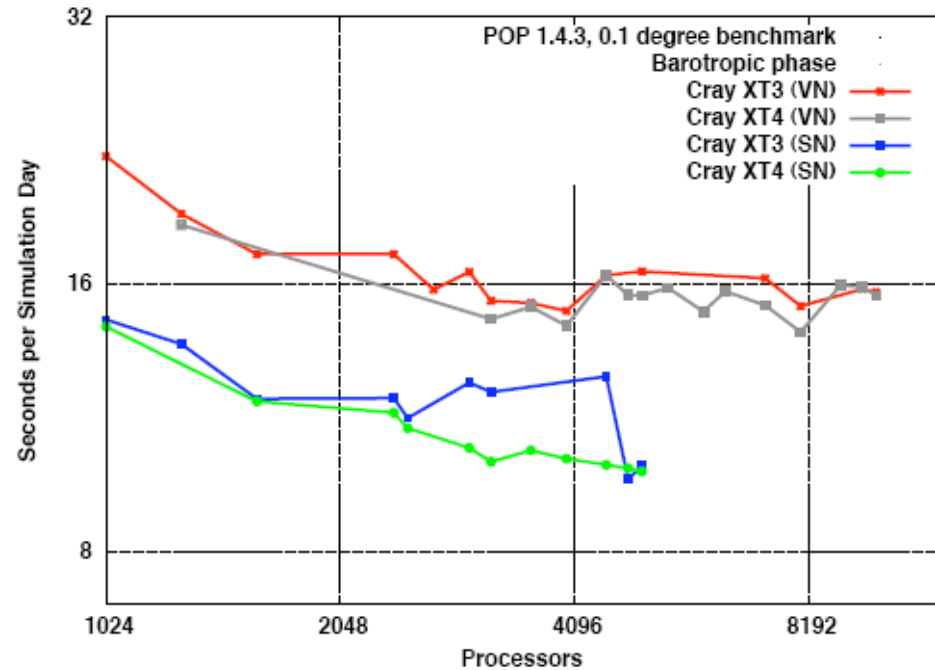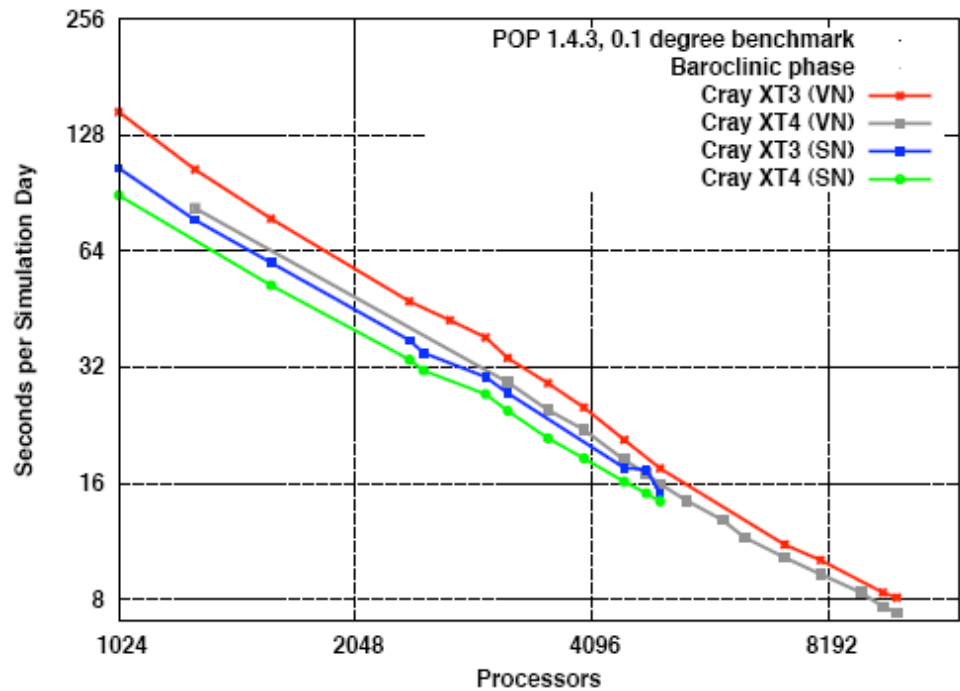
# User view of I/O

# Performance

# Performance > Matrix Multiply

# Performance > POP

# Performance > POP > Phases

# Performance Considerations

➡ Compiler options

➡ Page Size
- 4KB v. 2MB (def)

➡ SN v. VN

➡ Process Mapping
- Manage logical to physical placement of tasks
- Cray
  - MPICH_RANK_REORDER_METHOD
  - Wrap, smp-style

➡ Collectives
- MPI_COLL_OPT_ON

# Software

# Cray XT3/4 Software

➡ Operating system
 – Catamount
  • Lightweight OS
 – Compute Node Linux (in testing)
  • Derived from Linux
  • Targeting quad core release
  • More functionality
   – User threads
➡ Filesystem
 – Lustre
➡ Tools
 – Performance, debugging
  • Cray PAT and Apprentice
  • Tau
  • PAPI
  • MPIP
  • Totalview
➡ Compilers
 – PGI, Pathscale
➡ PBS/Moab batch scheduler

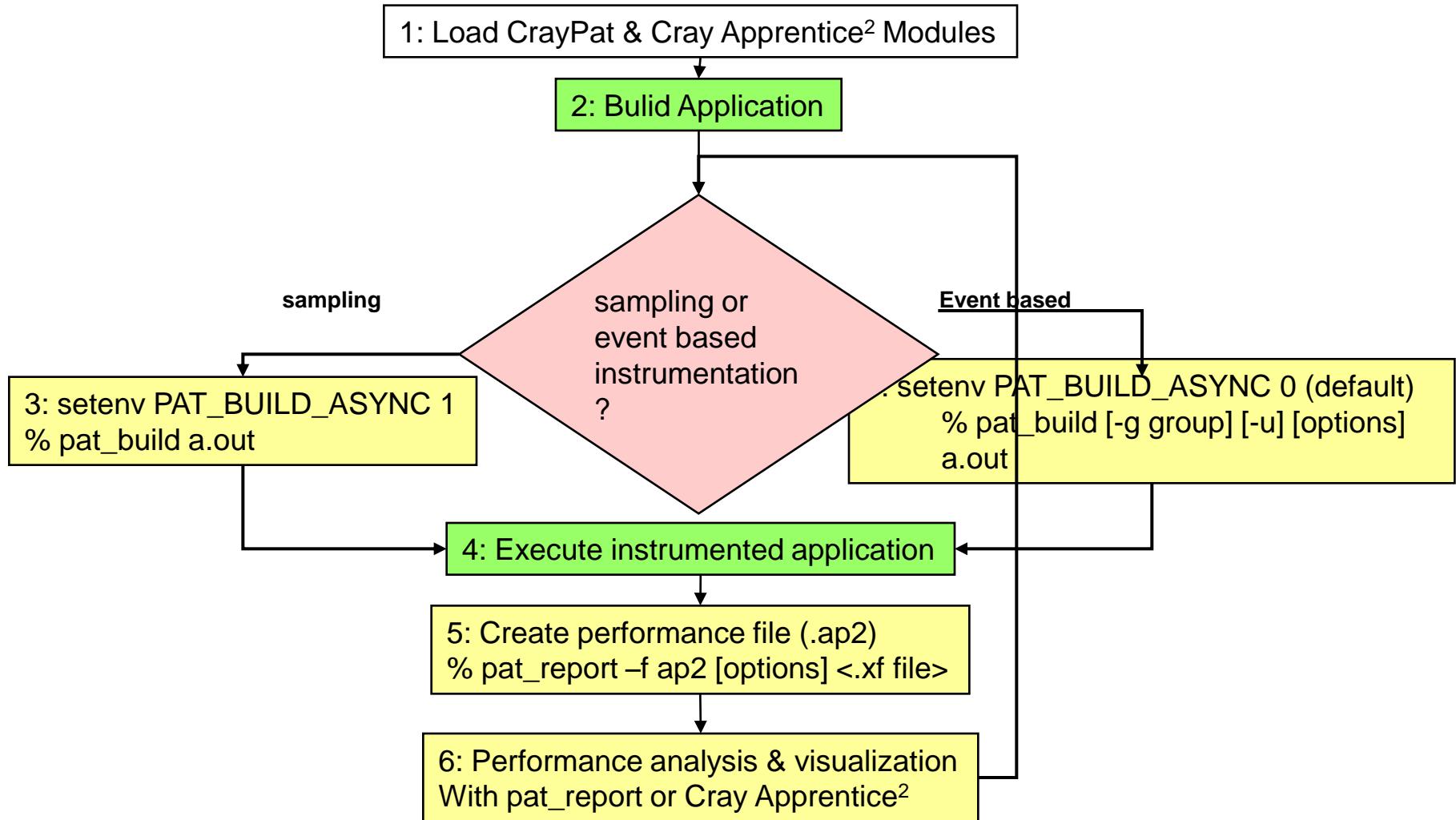# Cray Performance Analysis Infrastructure

➡ **CrayPat**

- pat_build: Utility for automatic application instrumentation
  - No source code modification required
- run-time library for measurements
  - transparent to the user
- pat_report:
  - Performance reports
  - Generation of file for performance visualization
- pat_help: Runtime help utility

➡ **Cray Apprentice[2]**

- Graphical performance analysis and visualization tool
  - Can be used off-line on Linux system

Cray Inc.

# Performance Analysis with CrayPat & Cray Apprentice[2]

1: Load CrayPat & Cray Apprentice[2] Modules

2: Bulid Application

sampling · sampling or event based instrumentation ? · **Event based**

3: setenv PAT_BUILD_ASYNC 1
% pat_build a.out

: setenv PAT_BUILD_ASYNC 0 (default)
% pat_build [-g group] [-u] [options]
a.out

4: Execute instrumented application

5: Create performance file (.ap2)
% pat_report –f ap2 [options] <.xf file>

6: Performance analysis & visualization
With pat_report or Cray Apprentice[2]

# Performance Metrics Available in pat_report

➡ **Profile by groups**
- – Threshold
- – Load imbalance information

➡ **Function Profile**
- – Flat profile
- – Call Tree view
- – Callers view
- – Hardware counters information

➡ **MPI Profiler**
- – MPI Load balance
- – MPI Stats by bin

➡ **I/O Statistics**
- – Read and Write Statistics

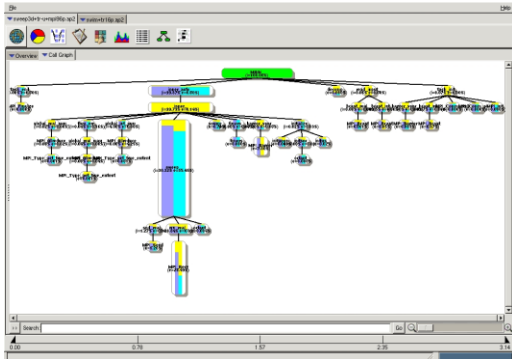➡ **Heap Statistics**
- – High water mark
- – Memory leaks

# CrayPat API

➡ CrayPat performs automatic instrumentation at function level

➡ The CrayPat API can be used for fine grain instrumentation

- Fortran
  - call PAT_region_begin(id, "label", ierr)
  - DO Work
  - call PAT_region_end(id, ierr)

- C
  - include <pat_api.h>
  - …
  - ierr = PAT_region_begin(id, "label");
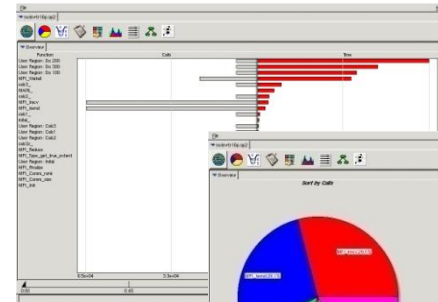  - DO_Work();
  - ierr = PAT_region_end(id);
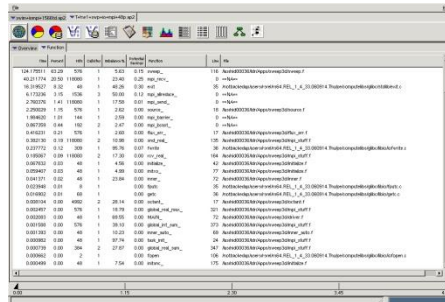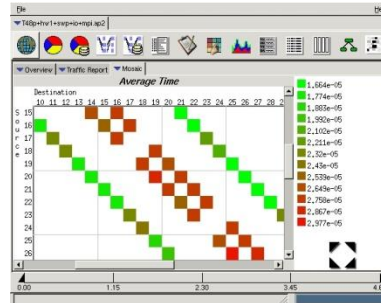
# Cray Apprentice²

**Function Profile**
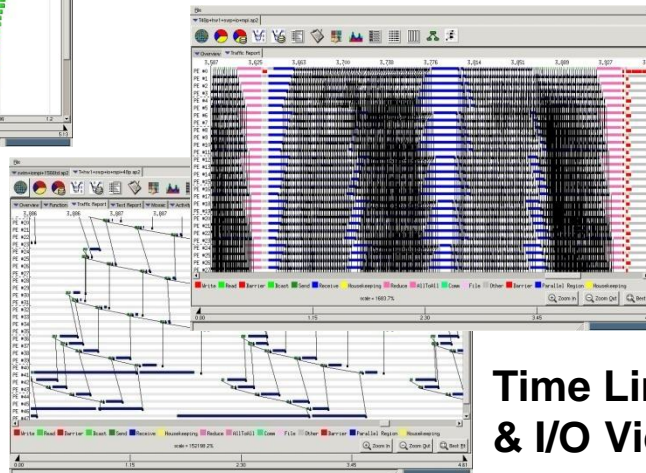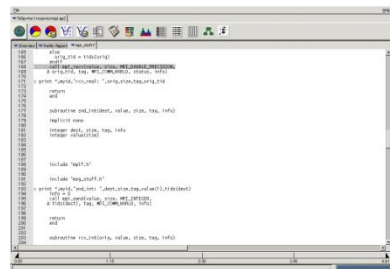
**Call Graph Profile**

**Function Overview**
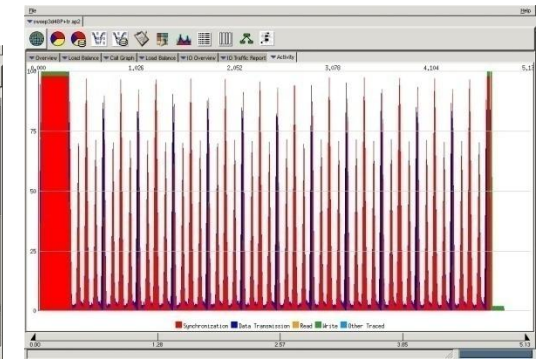
**Load balance views**

**Pair-wise Communication View**

**Source code mapping**

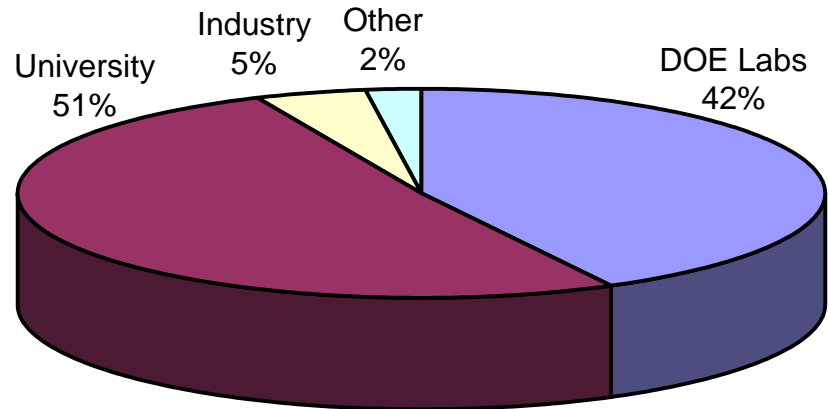**Time Line & I/O Views**

**Communication & I/O Activity View**

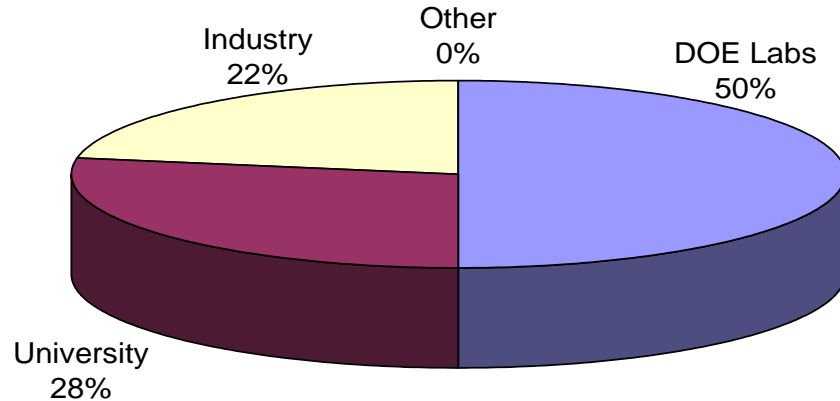# Usage

# NCCS Demographics in 2006

| 2006 Projects | | | |
|---|---|---|---|
| Accelerator physics | 1 | Engineering | 1 |
| Astrophysics | 3 | Fusion | 4 |
| Chemistry | 1 | High energy physics | 1 |
| Climate change | 3 | Biology | 2 |
| Combustion | 1 | Materials science | 2 |
| Computer science | 2 | Nuclear physics | 1 |

**LCF Users by Affiliation**



University 51%
Industry 5%
Other 2%
DOE Labs 42%

**FY 2006 Phoenix Usage by Affiliation**



Industry 22%
Other 0%
DOE Labs 50%
University 28%

**FY 2006 Jaguar Usage by Affiliation**



University 50%
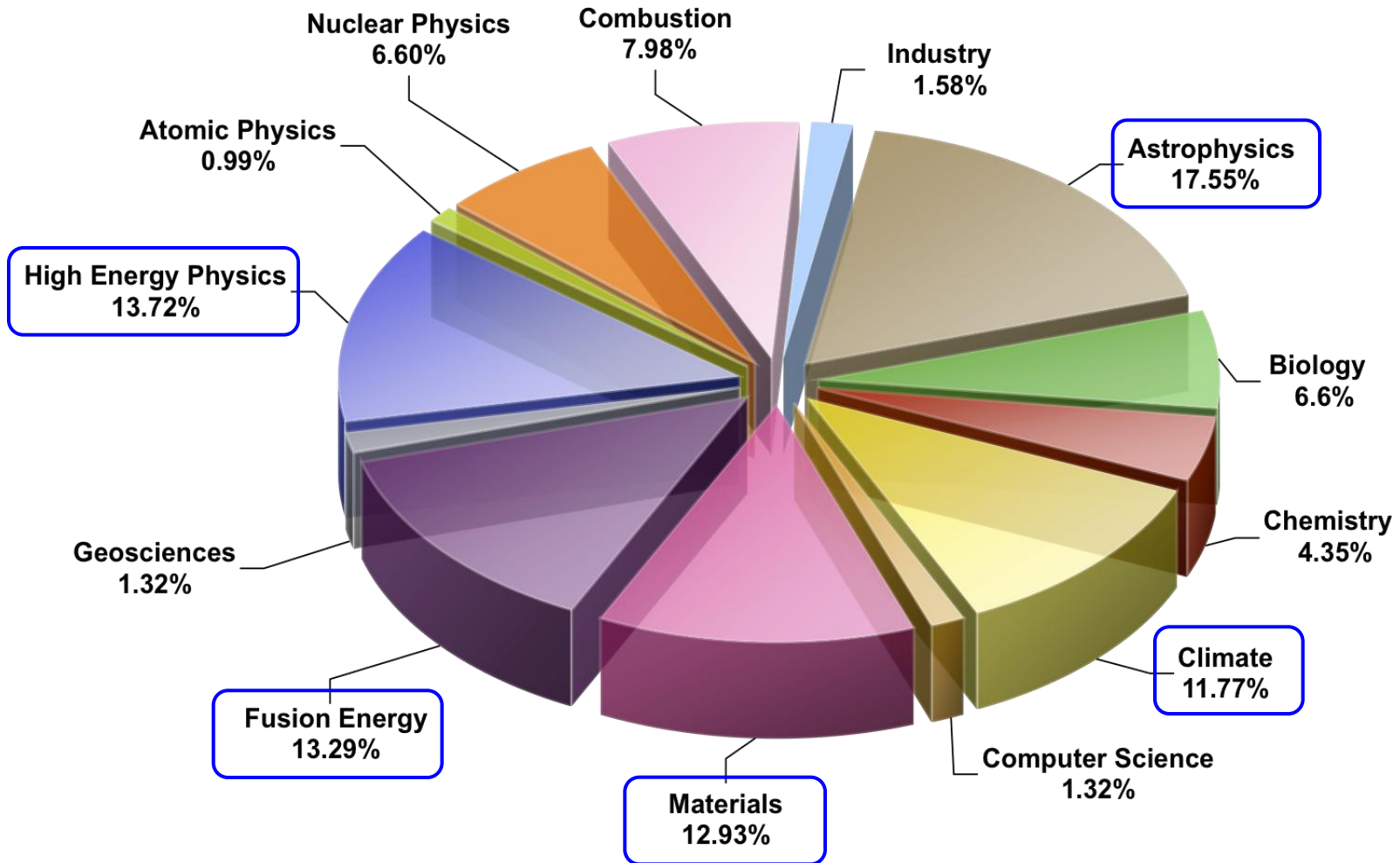Industry 0%
Other 0%
DOE Labs 50%

# Simulation Job Size Distribution for Science Applications on the ORNL Cray XT3 in 2006

# 2007 INCITE Allocations at the ORNL LCF
## Breakdown by Discipline



Pie chart: 2007 INCITE Allocations at the ORNL LCF — Breakdown by Discipline

- Nuclear Physics 6.60%
- Combustion 7.98%
- Industry 1.58%
- Atomic Physics 0.99%
- Astrophysics 17.55%
- High Energy Physics 13.72%
- Biology 6.6%
- Geosciences 1.32%
- Chemistry 4.35%
- Fusion Energy 13.29%
- Climate 11.77%
- Materials 12.93%
- Computer Science 1.32%

Circled domains make up 70% of total allocation

# NERSC Cray XT4

➡ Very similar to ORNL platform

➡ Main differences

    – XT4 nodes (homogeneous system)

    – GPFS filesystem