

Jaguar: The World's Most Powerful Computer

Buddy Bland
Cray Users Group 2009 Meeting
Atlanta, Georgia
May 5, 2009



U.S. DEPARTMENT OF
ENERGY

 **OAK RIDGE NATIONAL LABORATORY**
MANAGED BY UT-BATTELLE FOR THE DEPARTMENT OF ENERGY

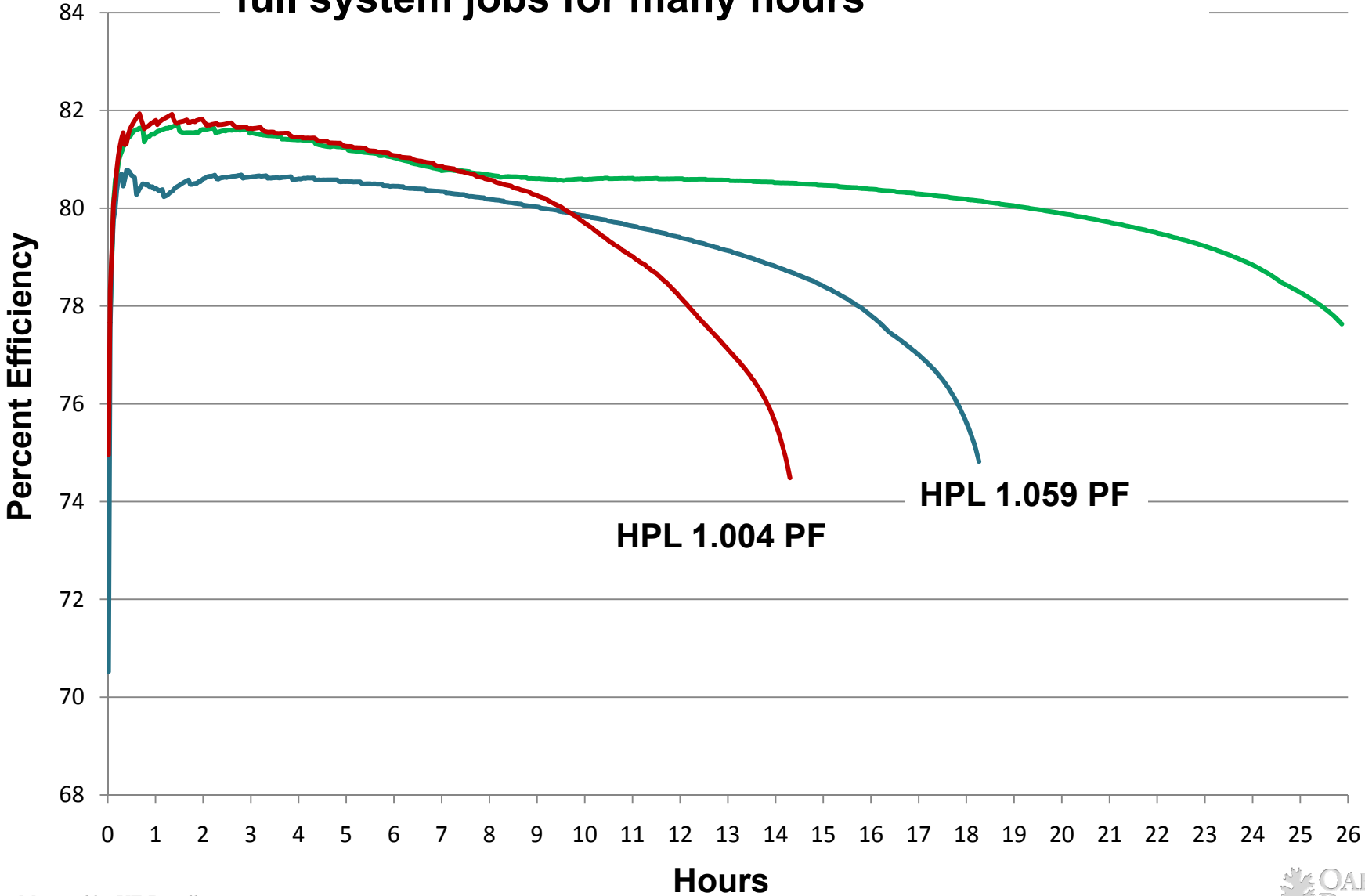
Outstanding launch for petascale computing in Office of Science and ORNL at SC'08

Only 41 days after assembly of a totally new 150,000 core system

- Jaguar beat the previous #1 performance on Top500 with the HPL benchmark running over 18 hours on the entire system
- Jaguar had two *real* applications running over 1 PF
 - DCA++ 1.35 PF Superconductivity problem
 - LSMS 1.05 PF Thermodynamics of magnetic nanoparticles problem

Cray XT5 “Jaguar” is showing impressive stability

Within days of delivery, the system was running full system jobs for many hours

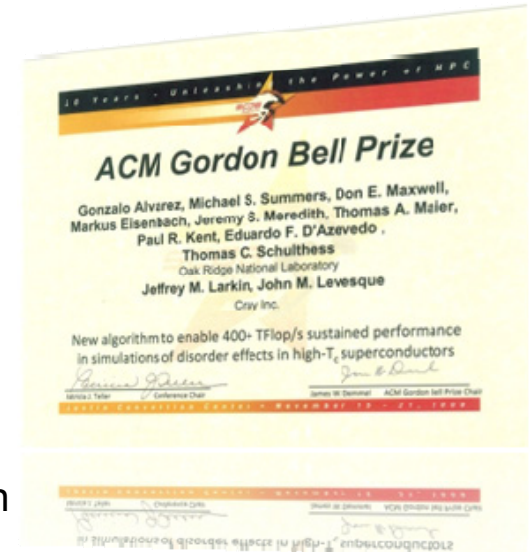


Gordon Bell prize awarded to ORNL team



Three of six GB finalist ran on Jaguar

- A team led by ORNL's Thomas Schulthess received the prestigious 2008 Association for Computing Machinery (ACM) Gordon Bell Prize at SC08
- For attaining fastest performance ever in a scientific supercomputing application
- Simulation of superconductors achieved 1.352 petaflops on ORNL's Cray XT Jaguar supercomputer
- By modifying the algorithms and software design of the DCA++ code, the team was able to boost its performance tenfold



DCA++
talk
Tuesday
1:00

Gordon Bell Finalists

- | | |
|-------------|------|
| ✓ DCA++ | ORNL |
| ✓ LS3DF | LBNL |
| ✓ SPECFEM3D | SDSC |
| • RHEA | TACC |
| • SPaSM | LANL |
| • VPIC | LANL |



HPC Challenge Awards




- HPC Challenge awards are given out annually at the Supercomputing conference
- Awards in four categories, result published for two others; tests many aspects of the computer's performance and balance
- Must submit results for all benchmarks to be considered
- Unfortunately, ORNL team only had two days on the machine to get the results. Got a better G-FFT number (5.804) the next day. ORNL submitted only baseline (unoptimized) results.

G-HPL (TF)		EP-Stream (GB/s)		G-FFT (TF)		G-Random Access (GUPS)		EP-DGEMM (TF)		PTRANS (GB/s)	
ORNL	902	ORNL	330	ANL	5.08	ANL	103	ORNL	1,257	SNL	4,994
LLNL	259	LLNL	160	SNL	2.87	LLNL	35.5	ANL	362	LLNL	4,666
ANL	191	ANL	130	ORNL	2.77↑	SNL	33.6	LLNL	162	LLNL	2,626

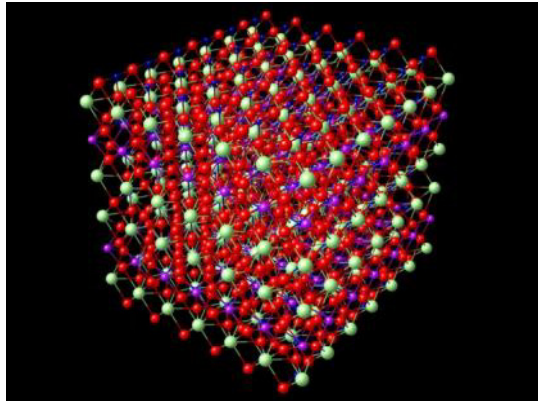
HPC CHALLENGE

Science Applications are Scaling on Jaguar

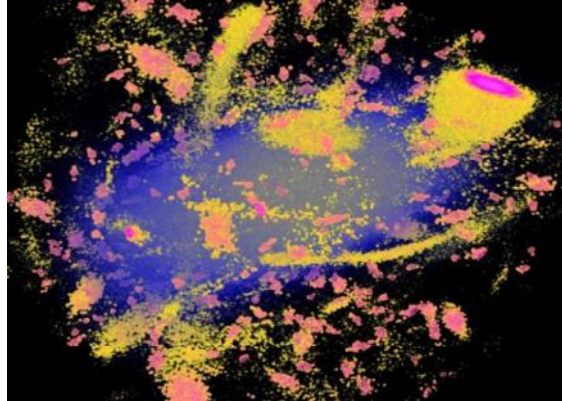
Science Area	Code	Contact	Cores	Total Performance	Notes
Materials	DCA++	Schulthess	150,144	1.3 PF*	Gordon Bell Winner 
Materials	LSMS	Eisenbach	149,580	1.05 PF	
Seismology	SPECFEM3D	Carrington	149,784	165 TF	Gordon Bell Finalist
Weather	WRF	Michalakes	150,000	50 TF	
Climate	POP	Jones	18,000	20 sim yrs/ day	
Combustion	S3D	Chen	144,000	83 TF	
Fusion	GTC	PPPL	102,000	20 billion Particles / sec	
Materials	LS3DF	Lin-Wang Wang	147,456	442 TF	Gordon Bell Winner 
Chemistry	NWChem	Apra	96,000	480 TF	
Chemistry	MADNESS	Harrison	140,000	550+ TF	

Enabling breakthrough science

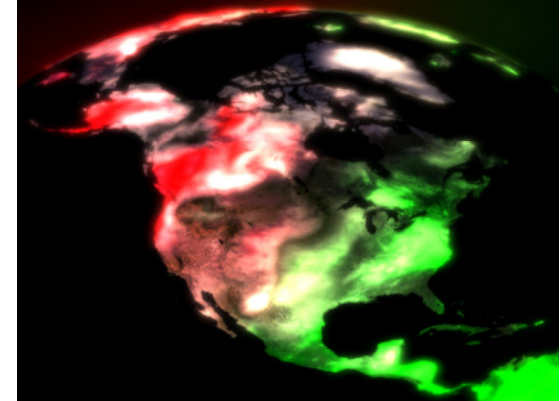
5 of top 10 ASCR science accomplishments in the past 18 months used LCF resources and staff



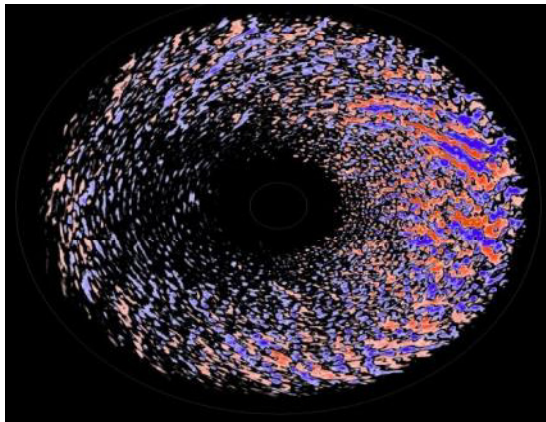
Electron pairing in HTSC cuprates
PRL (2007, 2008)



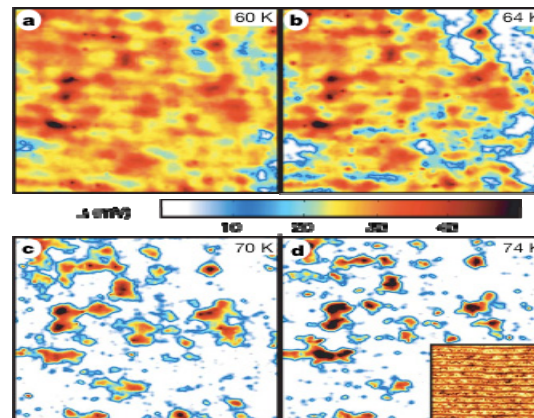
Shining a light on dark matter
Nature **454**, 735 (2008)



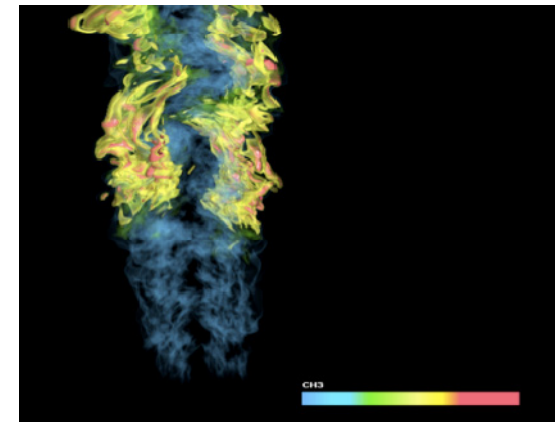
Modeling the full earth system



Fusion: Taming turbulent heat loss
PRL **99**, *Phys. Plasmas* **14**



Nanoscale nonhomogeneities in high-temperature superconductors
Winner of Gordon Bell prize

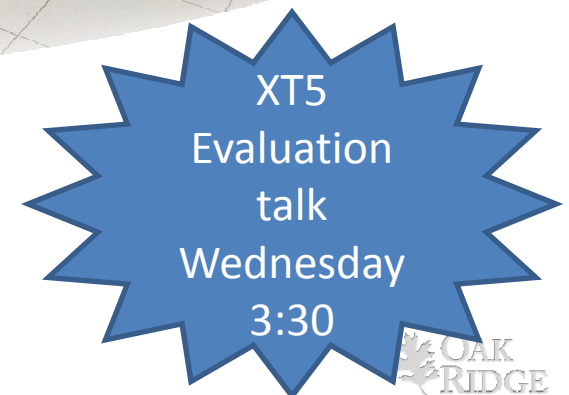


Stabilizing a lifted flame
Combust. Flame (2008)

Jaguar: World's most powerful computer Designed for science from the ground up



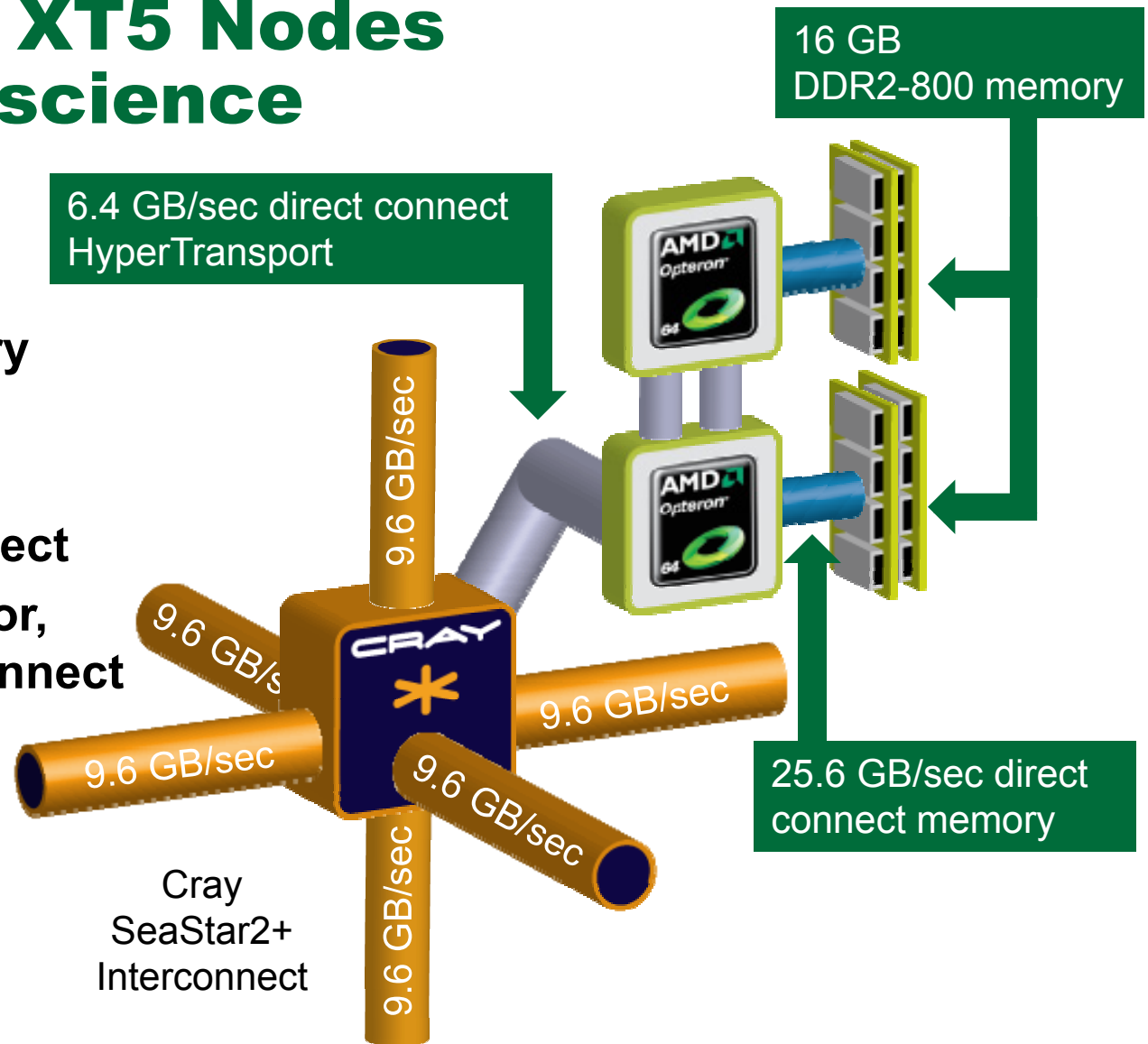
Peak performance	1.645 petaflops
System memory	362 terabytes
Disk space	10.7 petabytes
Disk bandwidth	200+ gigabytes/second



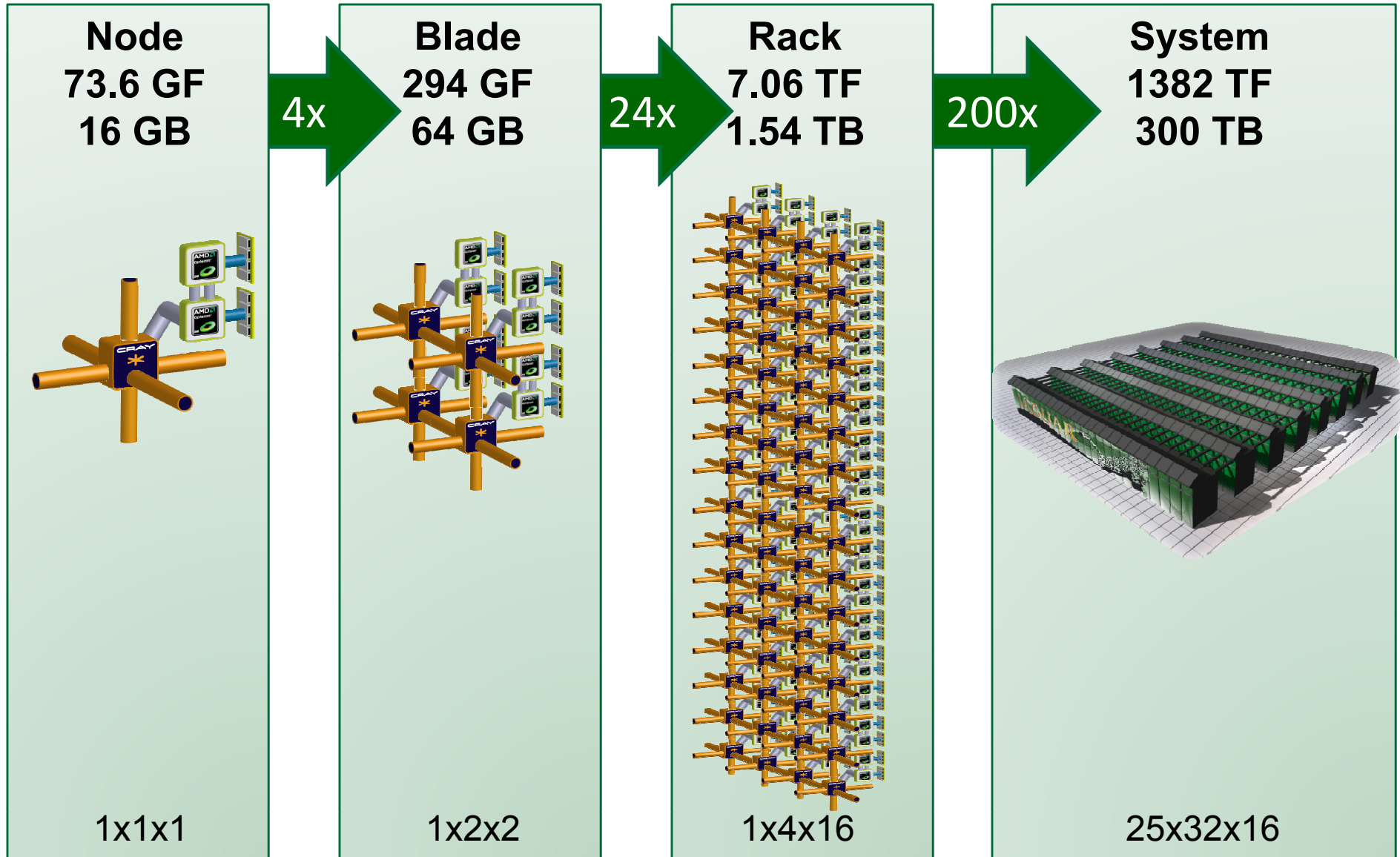
Jaguar's Cray XT5 Nodes Designed for science

- Powerful node improves scalability
- Large shared memory
- OpenMP Support
- Low latency, High bandwidth interconnect
- Upgradable processor, memory, and interconnect

GFLOPS	76.3
Memory (GB)	16
Cores	8
SeaStar2+	1



Building the Cray XT5 System



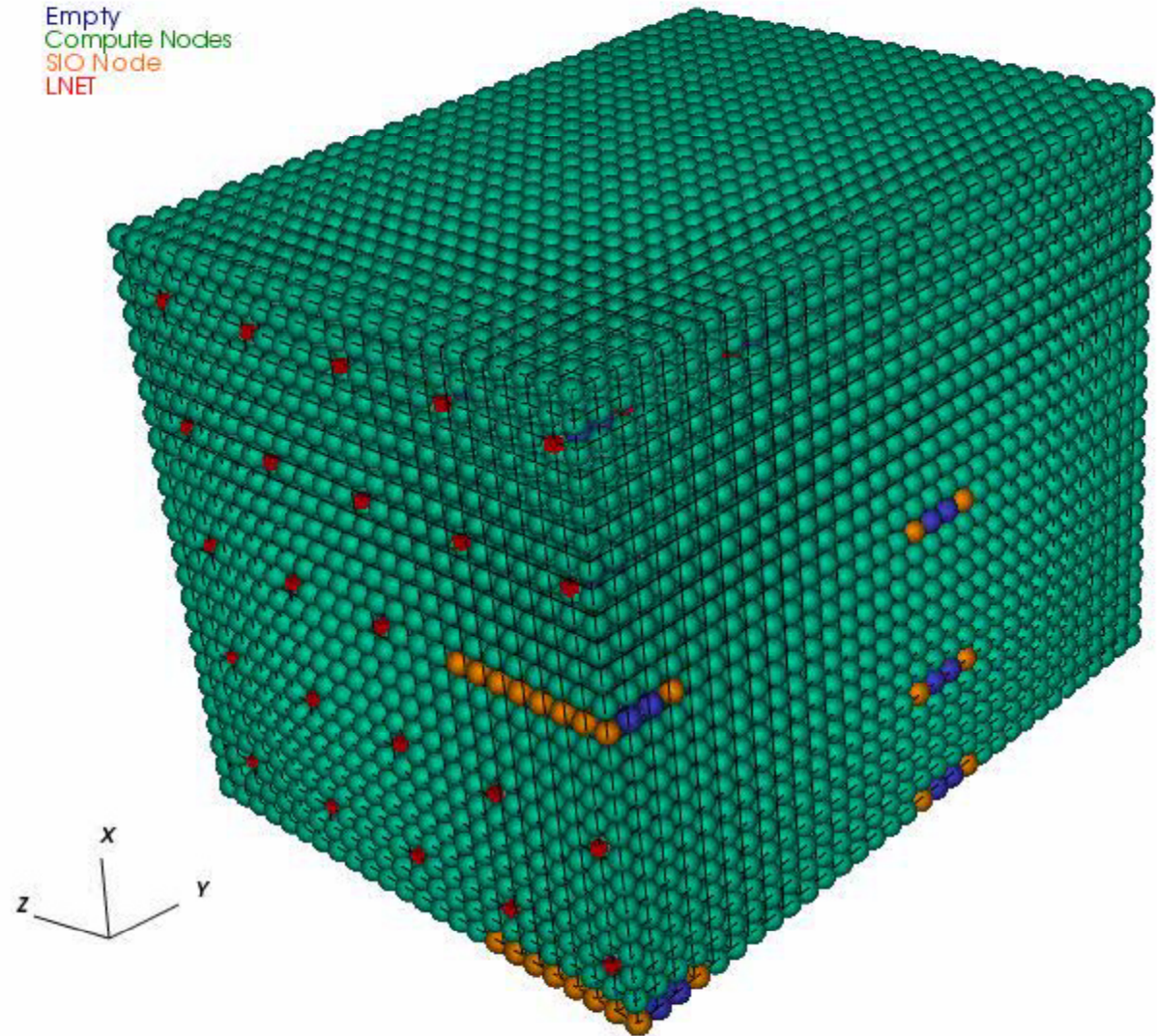
XT5 I/O Configuration Driven by application needs

XT5 Topology

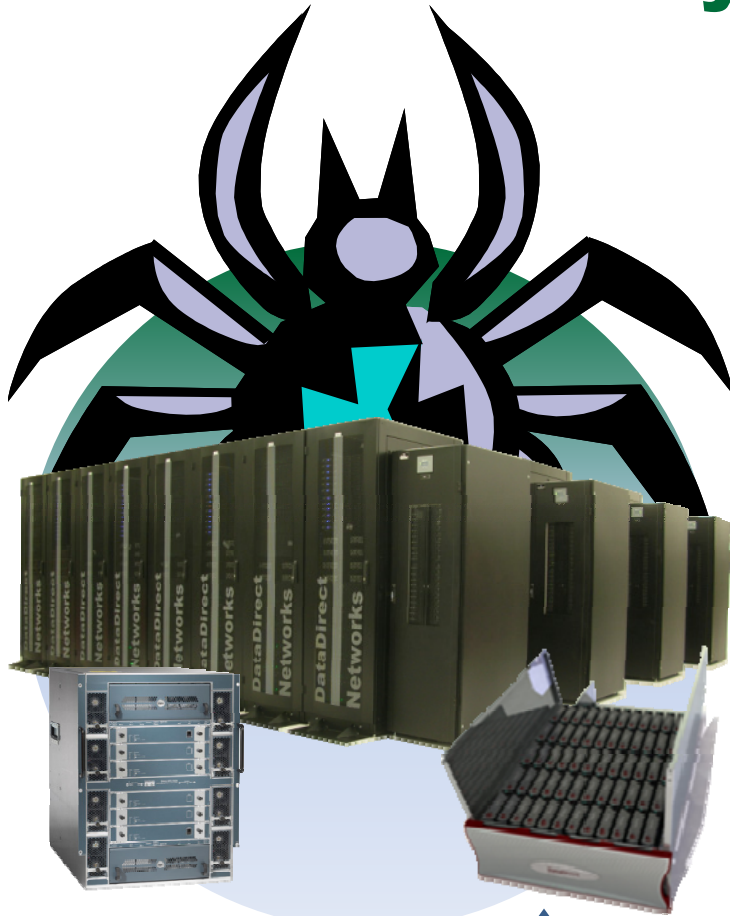
Features of I/O nodes

- 192 I/O nodes
- Each connected via non-blocking 4x DDR Infiniband to Lustre Object Storage Servers
- Fabric connections provides redundant paths
- Each OSS provide 1.25 GB/s
- I/O nodes spread throughout the 3-D torus to prevent hot-spots

Empty
Compute Nodes
I/O Node
LNET



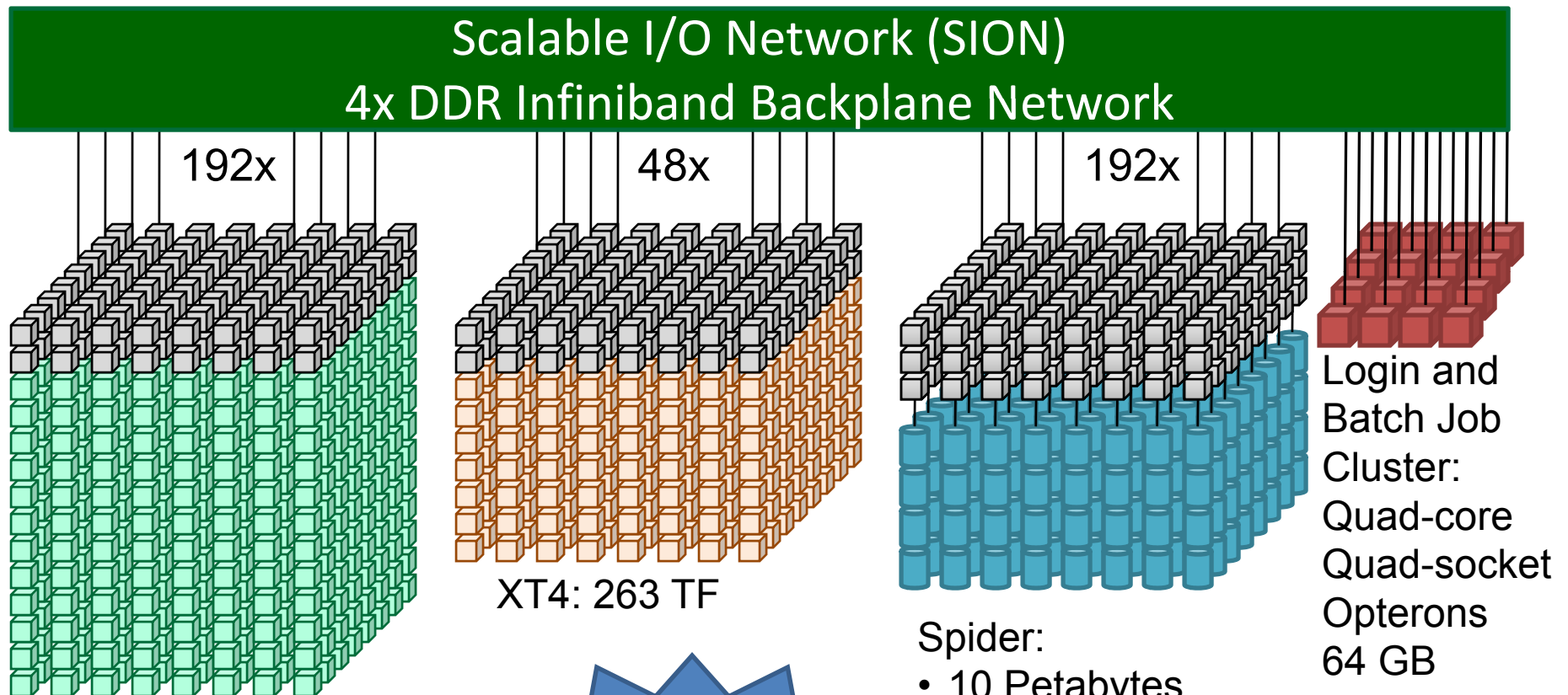
Center-wide File System



Spider talk
was
Monday at
2:00

- “Spider” provides a shared, parallel file system for all systems
 - Based on Lustre file system
- Demonstrated bandwidth of over 200 GB/s
- Over 10 PB of RAID-6 Capacity
 - 13,440 1-TB SATA Drives
- 192 Storage servers
 - 3 TB of memory
- Available from all systems via our high-performance scalable I/O network
 - Over 3,000 InfiniBand ports
 - Over 3 miles of cables
 - Scales as storage grows
- Undergoing friendly user checkout with deployment expected in summer 2009

Combine the XT5, XT4, and Spider with a Login Cluster to complete Jaguar



XT5: 1,382 TF

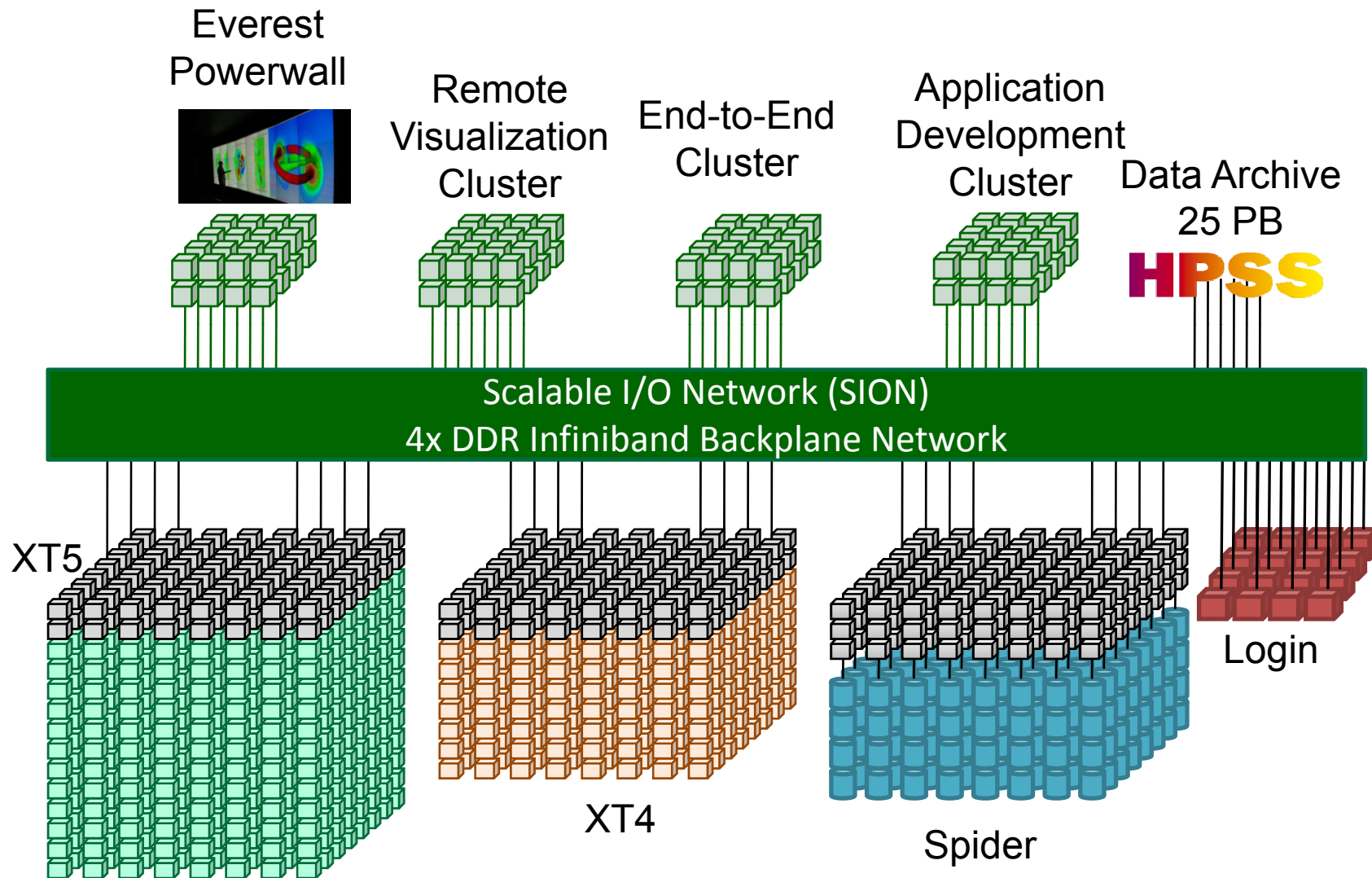
XT4: 263 TF

- Spider:
- 10 Petabytes
 - 192 OSS Nodes
 - 48 DDN 9900 Couplets
 - 13,440 disks

Login and Batch Job Cluster:
Quad-core
Quad-socket
Opteron
64 GB

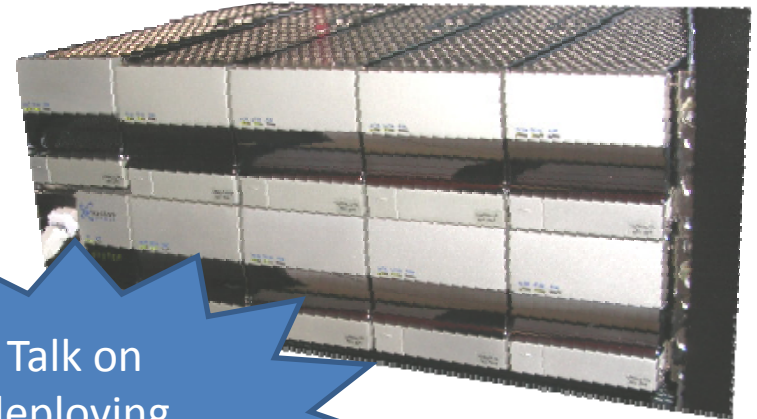
Talk on
integrating
XT4 and XT5
Thursday 8:30

Completing the Simulation Environment to meet the science requirements



XT5 Innovations: 480 volt power to the cabinet

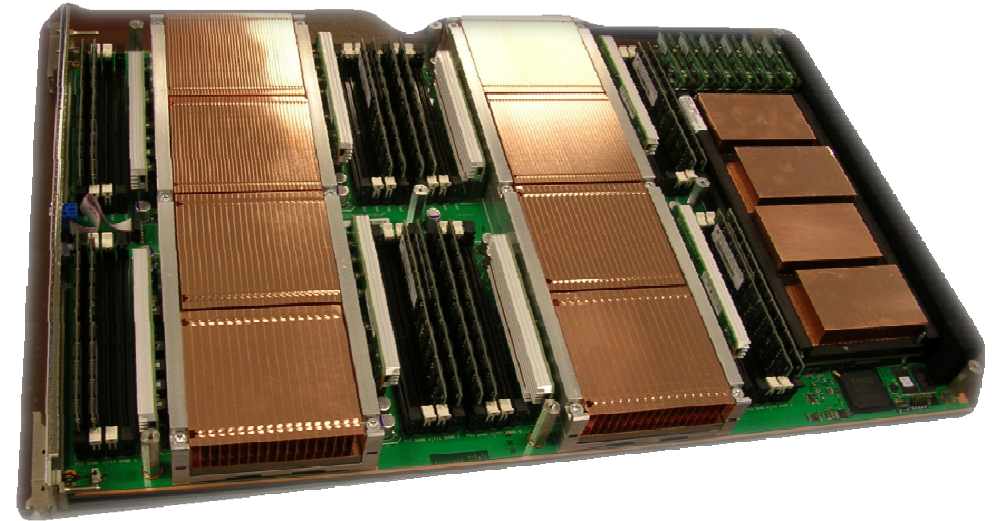
- Saved about \$1M in site prep costs in copper and circuit breakers
- Saves in ongoing electrical power costs by reducing losses in transformers and wires
- Allows higher density cabinets which shrinks system size



Talk on
deploying
Jaguar
Thursday 9:30

High-density blades

- Eight Opteron Sockets
- 32 DIMM slots
- 4 SeaStar2+ interconnect chips
- Variable pitch heat sinks



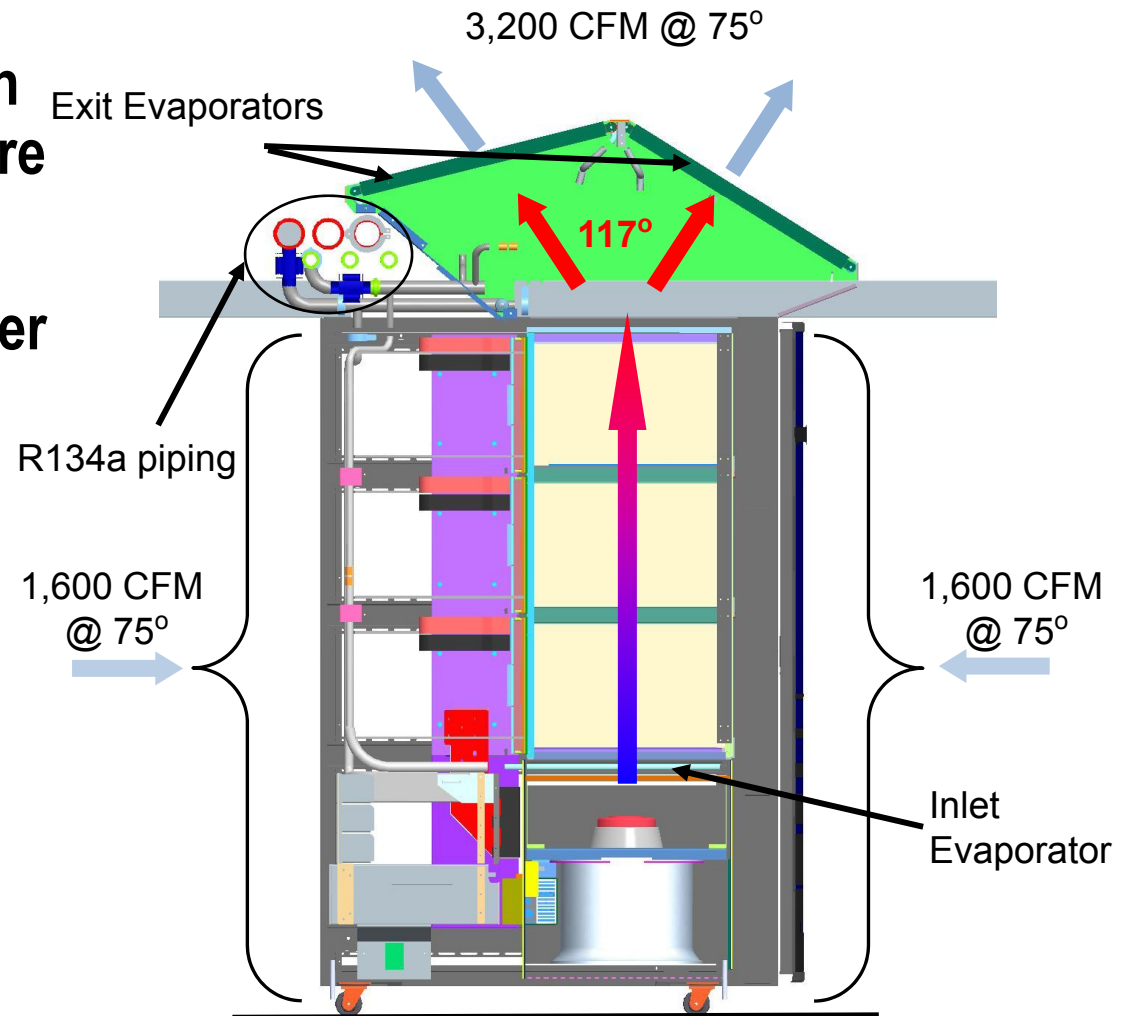
Single high-reliability fan

- Higher reliability than separate muffin-fans on each blade
- Custom designed turbine for high air-flow
- Variable speed to save power



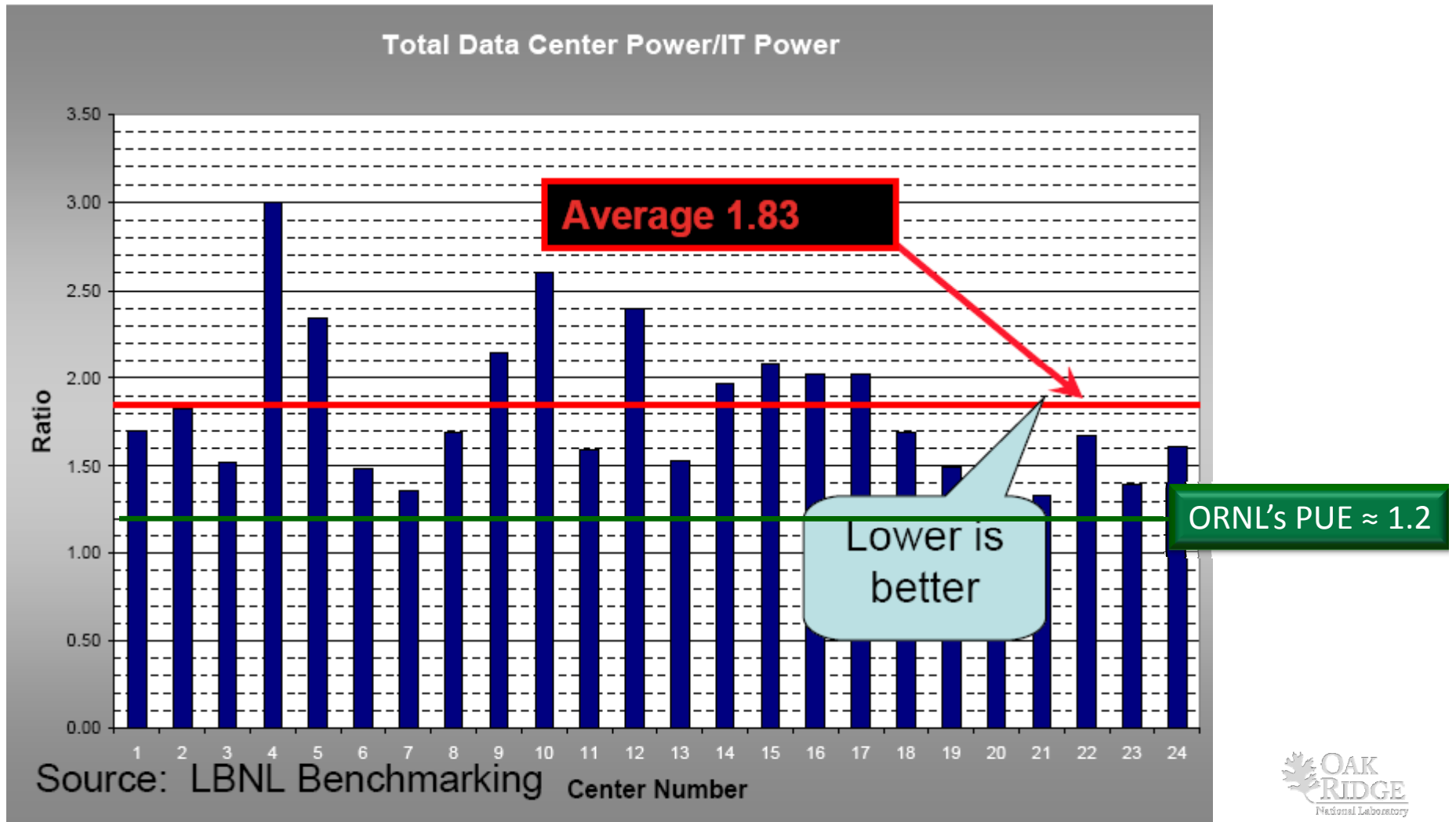
High Efficiency Liquid Cooling *Required* to build such a large system

- Newer Liquid Cooled design removes heat to liquid before it leaves the cabinet
- Saves about 900KW of power just in air movement and 2,500 ft² of floor space
- Phase change liquid to gas removes heat much more efficiently than water or air
- Each XDP heat exchanger replaces 2.5 CRAC units using one-tenth the power and floor space



Today, ORNL's facility is among the most efficient data centers

**Power Utilization Efficiency (PUE) =
Data Center power / IT equipment**



Electrical Systems Designed for efficiency

13,800 volt power into the building saves
on transmission losses



480 volt power to cabinets saves \$1M in
installation costs



High efficiency power supplies in the
cabinets



Flywheel based UPS for highest efficiency



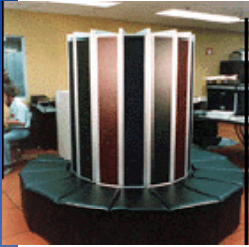
Just a word about power consumption

- The K-25 gaseous diffusion plant in Oak Ridge was completed in early 1945
- With 2 million square feet of space, it was the world's largest building at the time
- This entire building used 1,000 watts/ft². Most data centers today have much lower power density than this
- This one building used 2 gigawatts of power, which was about 10% of all the power generated in the U.S. at the time.



A bit of history about cooling and packaging

Power numbers in KW for a single CPU cabinet, not including SSD, IOS, HEU, or disks



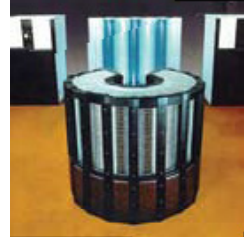
Cray-1

First Vector Supercomputer & first to utilize Freon cooling (150)



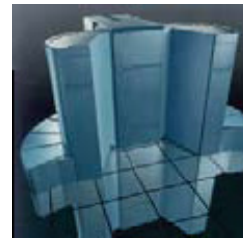
Cray X-MP

First vector multi-processor Supercomputer (160)



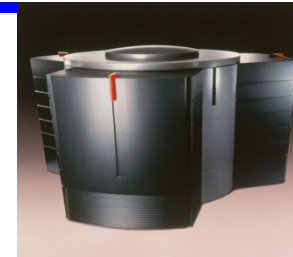
Cray-2

First Fluorinert Immersion cooled (200)



Cray Y-MP

First Supercomputer to sustain 1 GF, Fluorinert cold plates (145)



Cray C90

First Supercomputer with 1GF processor, Fluorinert cold plates (190)



Cray T90

First wireless supercomputer, Fluorinert immersion (345)



Cray T3E

First Supercomputer to sustain 1 TF, Fluorinert cold plate (45)



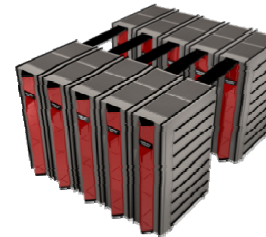
Cray X1/X1e

First Scalable Vector Supercomputer and first to utilize evaporative spray cooling (70)



Cray XT3/4

Highly scalable supercomputer, air cooled (20)



Cray XMT
First massively multithreaded supercomputer with extended memory semantics (25)



Cray XT5h

First Hybrid Supercomputer featuring scalable MPP, LC and Vector that utilized closed loop LC (45)



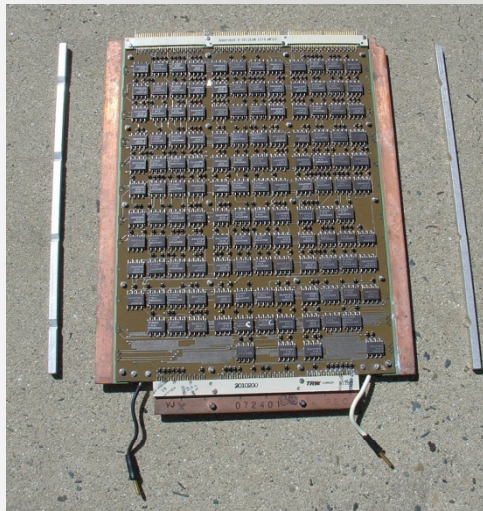
Cray XT5

First scalable system using R-134a cooling in top and bottom of the cabinet (40)



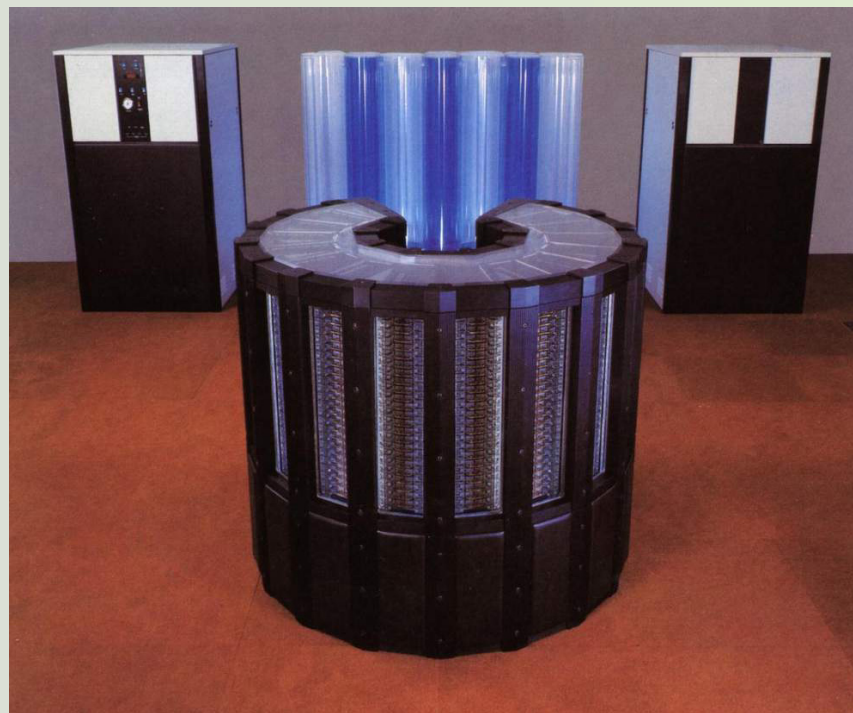
#1 Freon and Copper Cold Plates -1976

- Freon was used in conjunction with heat conducting plates
- Cray-1 and Cray XMP and I/O subsystems



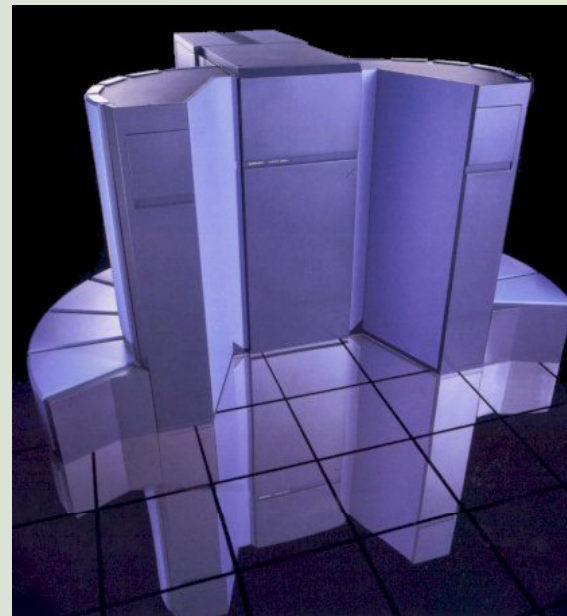
#2 Fluorinert Immersion -1986

- Initially used on the Cray-2 system
- Later used on the Cray T90 system and the Cray-3
- Entire computer is immersed in liquid
- Allowed tightly packed, 3 dimensional modules



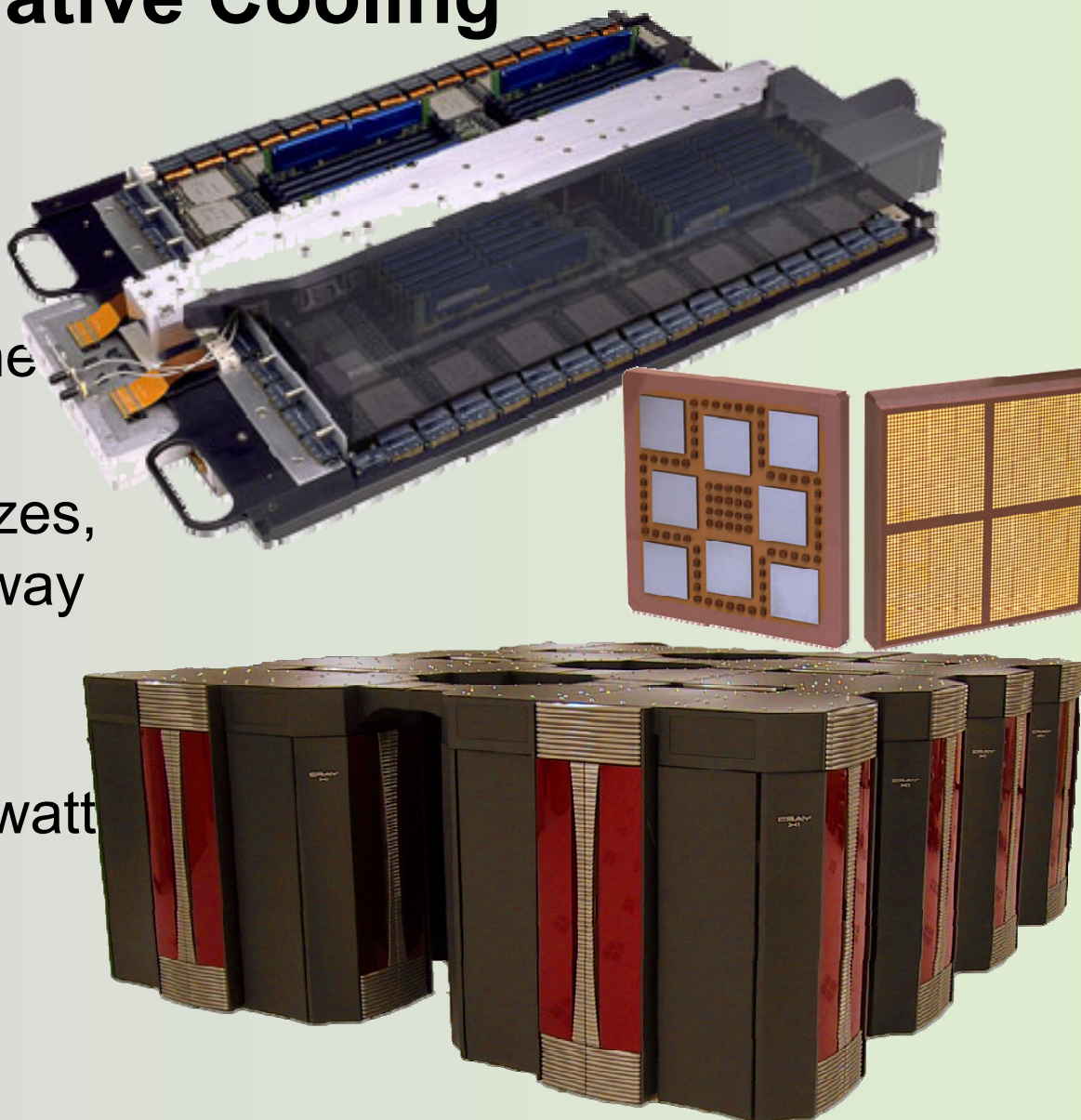
#3 Captive Fluorinert Cold Plates

- Used on the Cray Y-MP, Cray C90, Cray T3D and Cray T3E Systems
- Fluorinert circulated through a hollow cold-plate
- Fluorinert was used to minimize the chances of damage to components when the snap fittings were disconnected for servicing modules



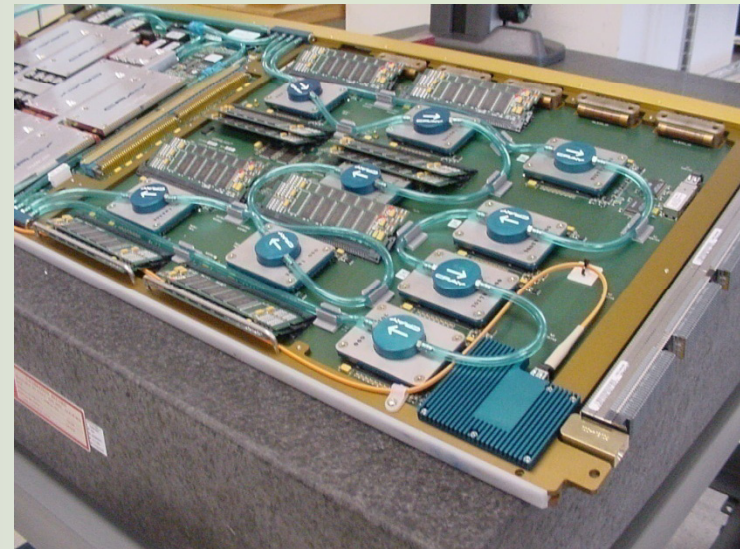
#4 Spray Evaporative Cooling

- Used on the Cray X1 processors
- A mist of Fluorinert is sprayed directly on the die
- The Fluorinert vaporizes, and heat is carried away via the latent heat of vaporization
- Used to cool a ~400 watt MCM



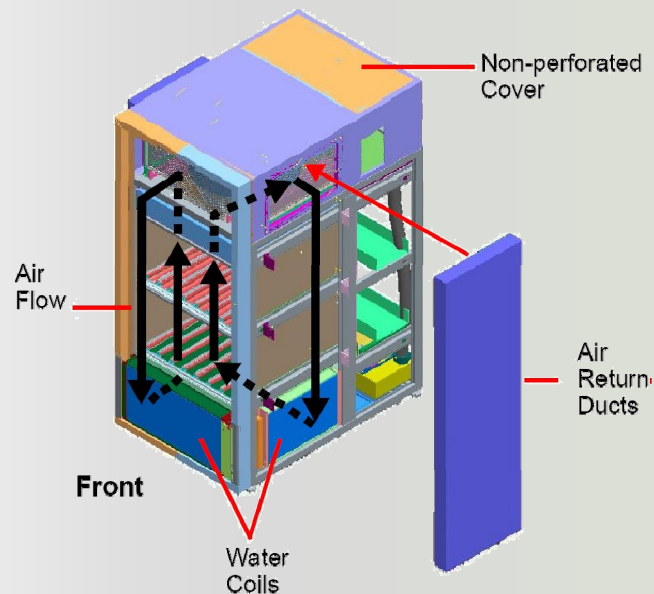
#5 Water Cap Cooling

- A water-filled heat-sink is mounted directly on an ASIC
- Used on the Cray MTA-2
- Designed to cool the custom ASICs in the machine
- Originally ran with water
- Later changed to Fluorinert because of organic growth in the fluid (and electrical problems induced by water flowing over dissimilar metals)



#6 Water Cooled Radiator

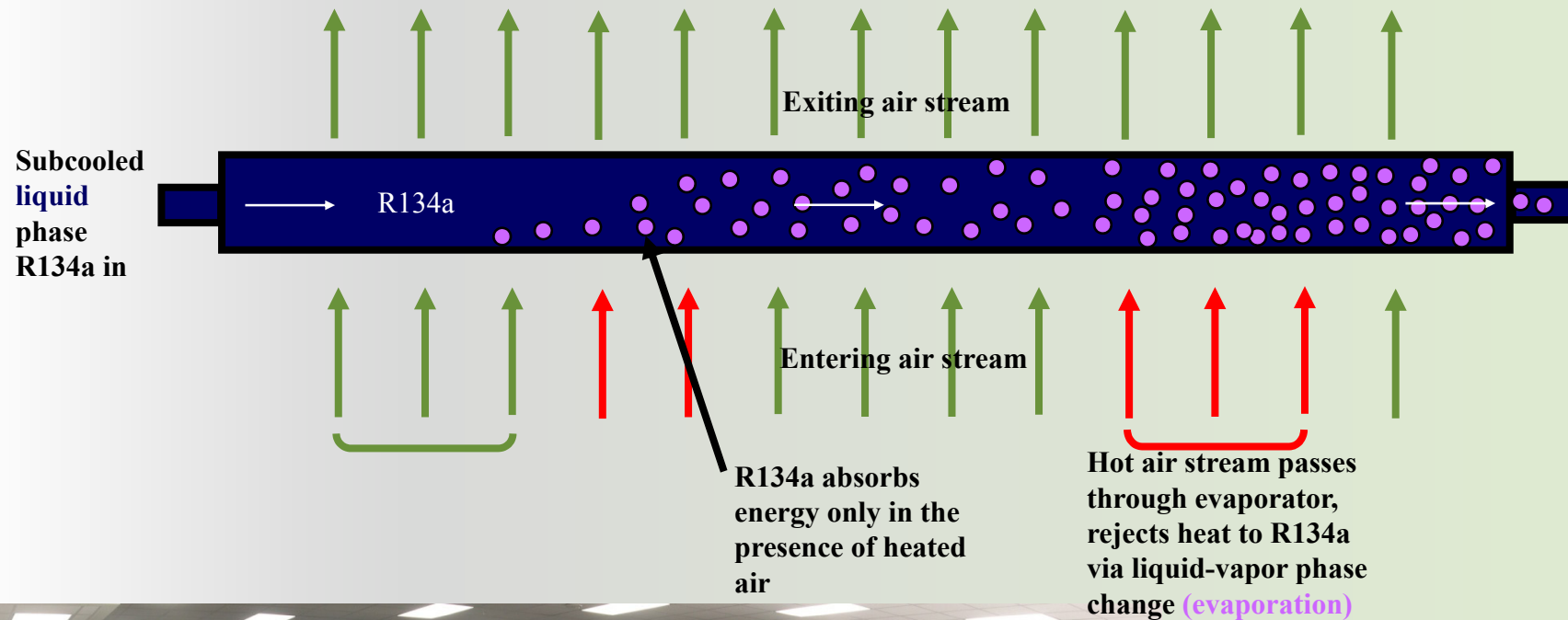
- Option on the Cray X2 vector processor cabinets
- Removes approximately 80% of the heat through chilled water
- Air is internally recirculated



Cooling Method #7

R134A Phase Change Evaporative Cooling

- Available on Cray XT5



Over 10x more effective than a water coil of similar size (phase change much more effective method to remove heat)

INCITE April 15th call for proposals

Call for large-scale, computationally intensive, high-impact research proposals

In 2010, powerful, leadership-class computing systems at DOE's Argonne National Laboratory and Oak Ridge National Laboratory will provide **over one billion** processor hours to a limited number of researchers nationwide.

The call is open to scientific researchers and research organizations, including industry; DOE Sponsorship is not required. **Deadline July 1st.**

INCITE awards help advance the state-of-the-art in areas such as

- Accelerator physics
- Astrophysics
- Chemical sciences
- Climate research
- Computer science
- Engineering
- Physics
- Environmental science
- Fusion energy
- Life sciences
- Materials science
- Nuclear physics, and more

For details about the DOE leadership computing facilities, see www.alcf.anl.gov and www.nccs.gov or contact INCITE@DOEleadershipcomputing.org to be added to an announcement distribution list.

Questions?

“We finally have a true leadership computer that enables us to run calculations impossible anywhere else in the world. The huge memory and raw compute power of Jaguar combine to transform the scale of computational chemistry.

Now that we have NWChem and MADNESS running robustly at the petascale, we are unleashing a flood of chemistry calculations that will produce insights into energy storage, catalysis, and functionalized nano-scale systems.”



*Robert Harrison
ORNL and University of Tennessee*