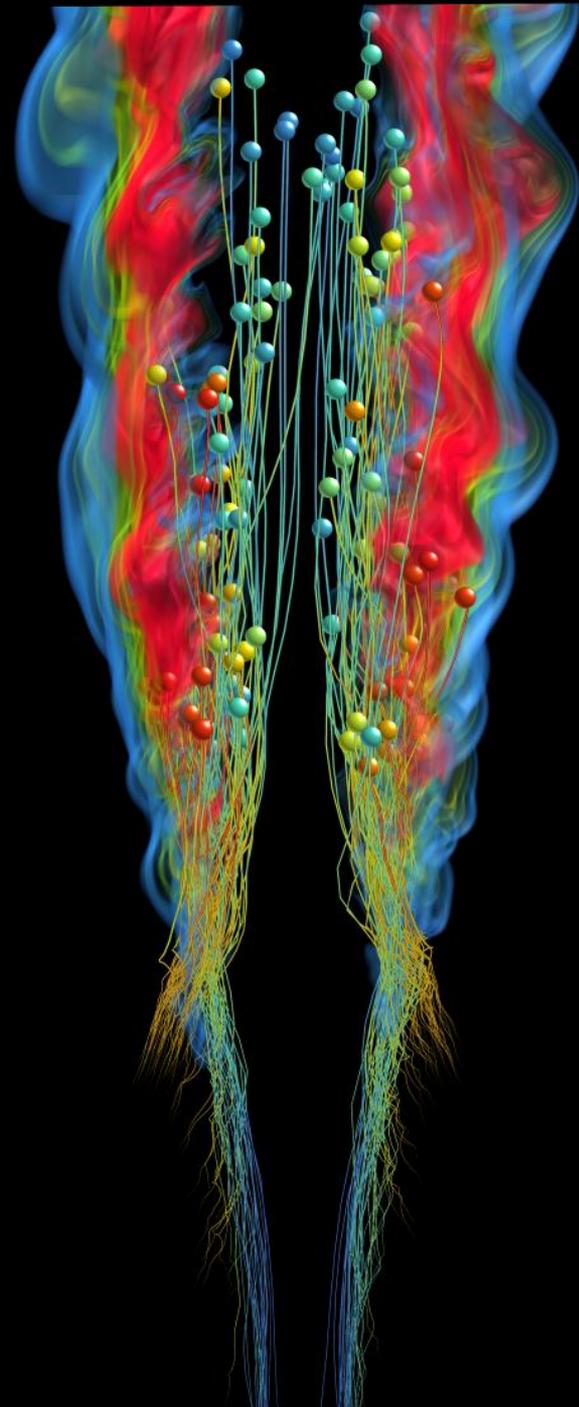# Keeneland: Toward Exascale Computational Science with Graphics Processors

**Jeffrey Vetter**

Presented to

**SC09**

**NVIDIA**

OAK RIDGE National Laboratory

Georgia Tech | College of Computing
Computational Science and Engineering

# Contributors

- **Jeffrey Vetter, Jack Dongarra, Richard Fujimoto, Thomas Schulthess, Karsten Schwan**, Sudha Yalamanchili, Kathlyn Boudwin, Jim Ferguson, Doug Hudson, Patricia Kovatch, Bruce Loftis, Jeremy Meredith, Jim Rogers, Philip Roth, Arlene Washington, Phil Andrews, Mark Fahey, Don Reed, Tracy Rafferty, Ursula Henderson, Terry Moore, and many others
- NVIDIA
- HP

- Keeneland Sponsor: NSF
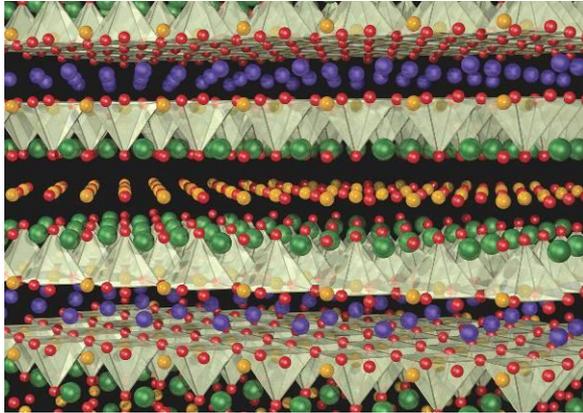- Other sponsors: DOE, DOD, DARPA

# Overview

- Predictive scientific simulation is important for scientific discovery
- HPC systems have been highly successful
- The HPC community has several (new) constraints
- Heterogeneous computing with GPUs offers some opportunities and challenges
- Newly awarded NSF partnership will provide heterogeneous supercomputing for open science
  - ORNL-Cray-NVIDIA announced system also using graphic processors
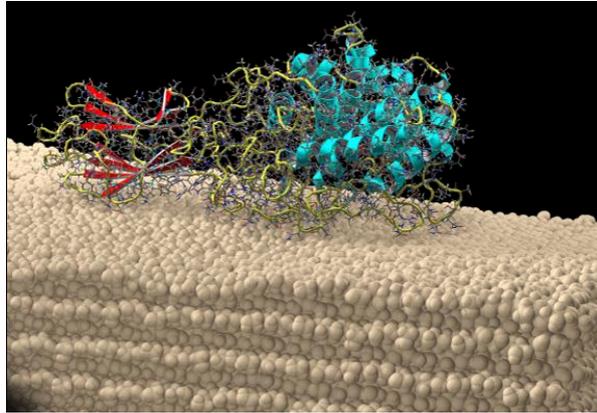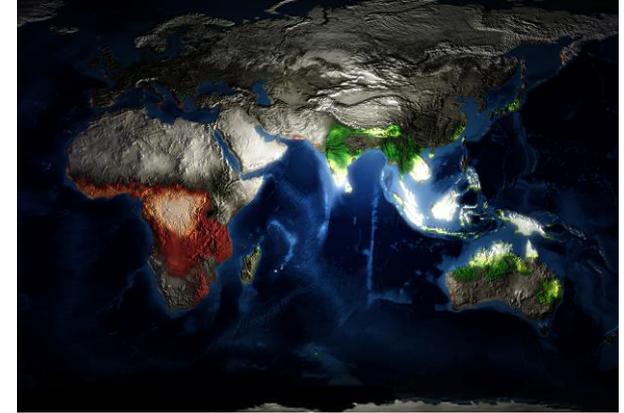
SCIENTIFIC SIMULATION

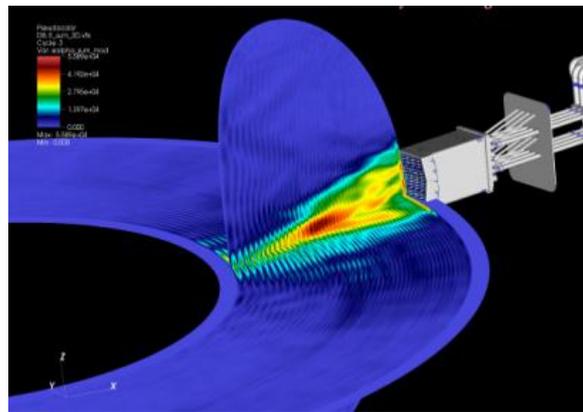# *Leadership computing is advancing scientific discovery*



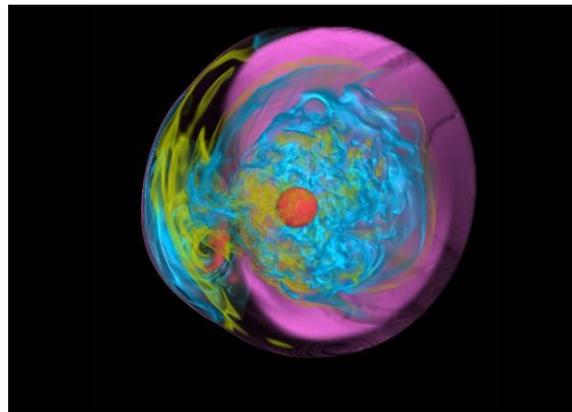**Resolved decades-long controversy about modeling physics of high temperature superconducting cuprates**



**New insights into protein structure and function leading to better understanding of cellulose-to-ethanol conversion**
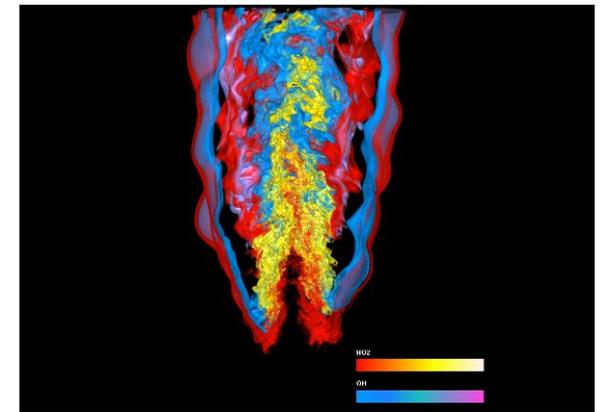


**Addition of vegetation models in climate code for global, dynamic $CO_2$ exploration**



**First fully 3D plasma simulations shed new light on engineering superheated ionic gas in ITER**



**Fundamental instability of supernova shocks discovered directly through simulation**
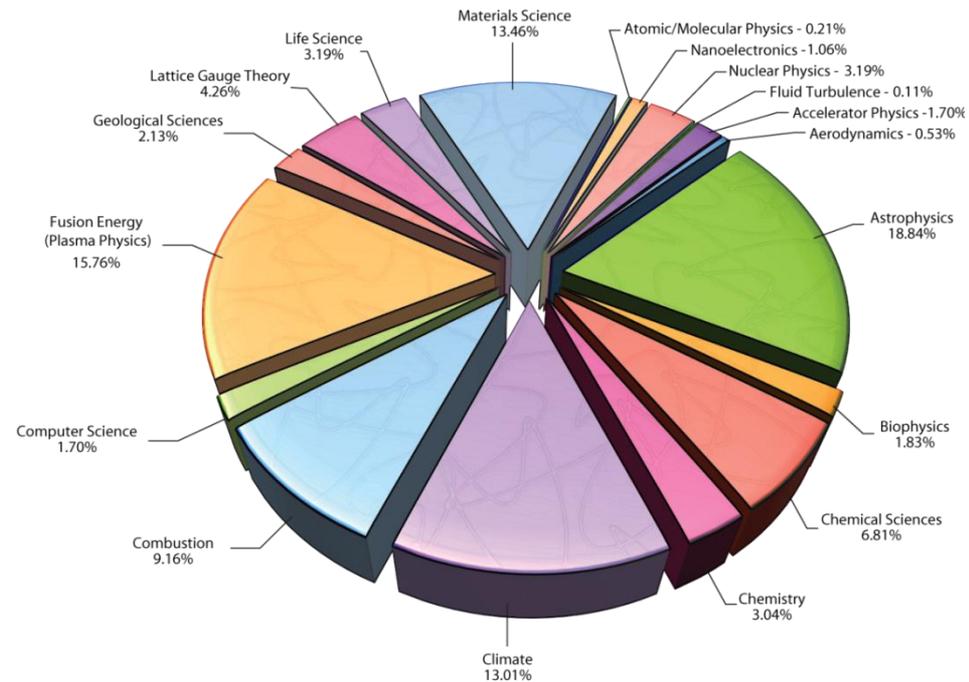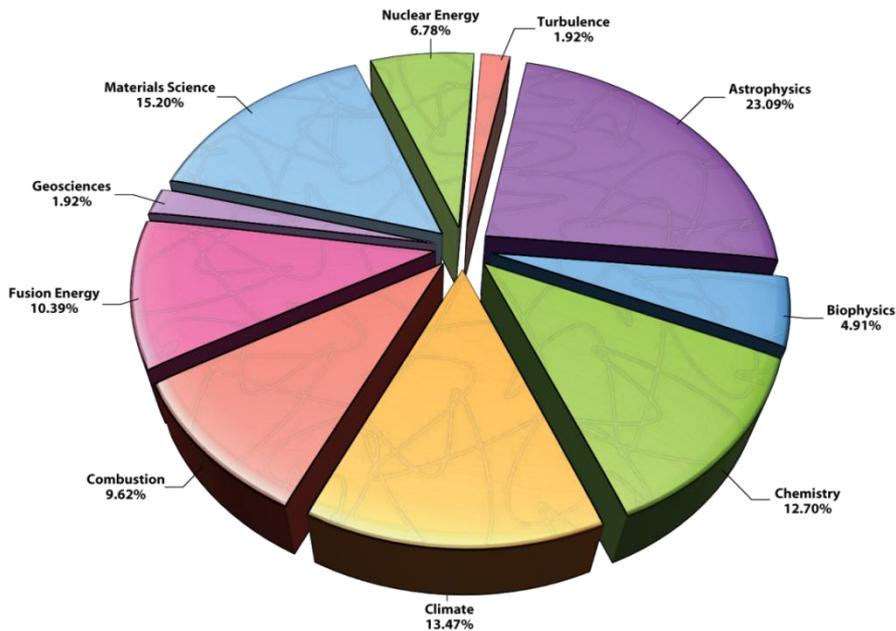


**First 3-D simulation of flame that resolves chemical composition, temperature, and flow**

# Jaguar XT5 science workload

- **Current: Transition to operations**
- **27 early-access projects, January–July**
  - Total allocation: 540.8M hours
  - Usage since January: 247M hours (159 users)
- **ALCC project: FY09 Joule metric applications**
  - Total allocation: 25M hours
  - Usage since January: 17M hours (17 users)

- **Fall 2009: General availability**
- **38 INCITE projects transitioned to system**
  - Total allocation: 469M hours
  - Usage (XT4) since January: 91M hours (448 users)
- **Discretionary projects: Climate AR5 production**
  - Total allocation: 80M hours
  - 2 projects: NCAR/DOE and NOAA (~50 users)



Left pie chart:
- Astrophysics 23.09%
- Biophysics 4.91%
- Chemistry 12.70%
- Climate 13.47%
- Combustion 9.62%
- Fusion Energy 10.39%
- Geosciences 1.92%
- Materials Science 15.20%
- Nuclear Energy 6.78%
- Turbulence 1.92%

Right pie chart:
- Materials Science 13.46%
- Atomic/Molecular Physics - 0.21%
- Nanoelectronics - 1.06%
- Nuclear Physics - 3.19%
- Fluid Turbulence - 0.11%
- Accelerator Physics - 1.70%
- Aerodynamics - 0.53%
- Life Science 3.19%
- Lattice Gauge Theory 4.26%
- Geological Sciences 2.13%
- Fusion Energy (Plasma Physics) 15.76%
- Astrophysics 18.84%
- Biophysics 1.83%
- Chemical Sciences 6.81%
- Chemistry 3.04%
- Computer Science 1.70%
- Combustion 9.16%
- Climate 13.01%

# *Highly visible science output*

**Physical Review Letters:**
High temperature superconductivity

**Combustion and Flame:**
3D flame simulation

**Nature:**
Astrophysics

**Physics of Plasmas:**
ICRF heating in ITER

**Transformational science enabled by advanced scientific computing**

# *Examples: Climate Modeling*

- **Assess scientific basis, impact, adaptation, vulnerability, mitigation**
  - > Observation and simulation play critical roles

- ***Profound impact on worldwide socioeconomic policy, energy, etc***
  - > UN Conference - Sep 09
  - > Copenhagen Climate conference – Dec 09

- **Intergovernmental Panel On Climate Change**
  - > Sponsored by UNEP and WMO
  - > Thousands of scientists from around the world

Muir Glacier, Alaska, August 13, 1941, photo by W.O. Field

Muir Glacier, Alaska, August 31, 2004, photo by B.F. Molnia

COP15 COPENHAGEN

UNITED NATIONS CLIMATE CHANGE CONFERENCE DEC 7-DEC 18 2009

http://www.ipcc.ch          http://en.cop15.dk/

# *Application Complexity*

- **Today's applications are quickly becoming more complex than earlier applications**
  - > Exascale Town Hall report [1]
- **Increased fidelity**
- **Answering more complex questions**
  - > Scientists [2,3] working on model that combines
    - – Climate model
    - – Energy economics
    - – Population movements, demographics

New physics capabilities are continuously integrated into ALE3D

- Dates are approximate, and represent *initial* capability
- List does not include computer science advancements (e.g. new databases, parsers, objects, etc…)

| Initial Capability | Physics package |
|---|---|
| 1994 | ALE explicit hydro w/ advection and slide surfaces (ASC code) |
| 1995 | Thermal transport (Topaz) |
| 1996 | Chemical reactions |
| 1997 | Implicit time stepping |
| 1998 | Shell (structural) elements |
| 1999 | Species diffusion of chemicals |
| 2000 | Incompressible flow |
| 2001 | Compressible flow |
| 2002 | MS material model library |
| 2002 | Particle tracking |

| | |
|---|---|
| 2003 | Built-in mesh generator |
| 2004 | Dislocation dynamics |
| 2004 | Beam elements |
| 2005 | Multiple implicit element types |
| 2005 | Multiphase flow hydro |
| 2005 | 2D axisymmetric |
| 2005 | Detonation Shock Dynamics (DSD) |
| 2006 | Magneto Hydrodynamics (MHD) |
| 2007 | Adaptive Mesh Refinement (AMR) |
| 2008 | Multiple time-integration options |
| 2008 | Auto contact |
| 2009? | Embedded grids |
| 2009? | Discrete Element Methods (DEM) |
| 2010? | Fully coupled AMR |

Still work in progress   Future plans

S&T / Comp / WCI

11

Source: LLNL ALE3D Team

---

Applications design and implementation are already complex!!
Writing and optimizing code for each new architecture and programming model is impractical (and only going to happen w/ heroic efforts/funding.)

[1] H. Simon, T. Zacharia, and R. Stevens, Eds., Modeling and Simulation at the Exascale for Energy and the Environment, 2007.
[2] S.W. Hadley, D.J. Erickson et al., "Responses of energy use to climate change: A climate modeling study," *Geophysical Research Letters*, 33(17), 2006.
[3] D. Vergano, "Air conditioning fuels warming," in USA Today, 2006.

# HPC LANDSCAPE TODAY
# JAGUAR AT ORNL

# *HPC Landscape Today*

- **2.3 PF Peak**
- **1.7 PF announced yesterday**

| Jaguar Specifications | Total | XT5 | XT4 |
|---|---|---|---|
| Peak Teraflops | 1,645 | 1,382 | 263 |
| Quad-Core AMD Opterons | 45,376 | 37,544 | 7,832 |
| AMD Opteron Cores | 181,504 | 150,176 | 31,328 |
| Compute Nodes | 26,604 | 18,772 | 7,832 |
| Memory (TB) | 362 | 300 | 62 |
| Disk Bandwidth (GB/s) | 284 | 240 | 44 |
| Disk Space (TB) | 10,750 | 10,000 | 750 |
| Interconnect Bandwidth (TB/s) | 532 | 374 | 157 |
| Floor Space (feet$^2$) | 5,800 | 4,400 | 1,400 |
| Cooling Technology | | Liquid | Air |

**Cray XT5 Node Architecture**

Cray SeaStar2+ Interconnect

http://nccs.gov

# Science applications are scaling on Jaguar

| Science area | Code | Contact | Cores | Performance | |
|---|---|---|---|---|---|
| **Materials** | **DCA++** | **Schulthess** | **150,144** | **1.3 PF MP** |  |
| Materials | LSMS | Eisenbach | 149,580 | 1.05 PF | |
| **Seismology** | **SPECFEM3D** | **Carrington** | **149,784** | **165 TF** |  |
| Weather | WRF | Michalakes | 150,000 | 50 TF | |
| Climate | POP | Jones | 18,000 | 20 simulation years/day | |
| Combustion | S3D | Chen | 144,000 | 83 TF | |
| Fusion | GTC | PPPL | 102,000 | 20 billion particles/second | |
| **Materials** | **LS3DF** | **Lin-Wang Wang** | **147,456** | **442 TF** |  |
| Chemistry | NWChem | Apra | 96,000 | 480 TF | |
| Chemistry | MADNESS | Harrison | 140,000 | 550+ TF | |

More apps to come in this year's Gordon Bell contest

THE ROAD AHEAD

# *Representative Roadmap to Exascale from ORNL*



Increasing computation requirements and resources

1 EF

100 PF > 250 PF

20 PF > 40 PF

0.6 -> 1 PF Cray XT(NSF- 1)

1 -> 2 PF Cray (LCF-2)

170 TF Cray XT4 (NSF-0)

50 TF > 100 TF > 250 TF Cray XT4 (LCF-1)

18.5 TF Cray X1E (LCF- 0)

ORNL Multi-Agency Computer Facility 260,000 ft$^2$

ORNL Multipurpose Research Facility

ORNL Computational Sciences Building

| 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 |

# ... facilities and cost issue ???



**Increasing computation requirements and resources**

1 EF

100 PF > 250 PF

20 PF  >  40 PF

0.6 -> 1 PF Cray XT(NSF- 1)

1 -> 2 PF Cray (LCF-2)

170 TF Cray XT4 (NSF-0)

50 TF > 100 TF > 250 TF Cray XT4 (LCF-1)

18.5 TF Cray X1E (LCF- 0)

**Increasing power and facilities**

?

ORNL Multi-Agency Computer Facility 260,000 ft²

ORNL Multipurpose Research Facility

ORNL Computational Sciences Building

| 2006 | 2007 | 2008 | 2009 | 2010 | 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 |

# Reality: Facilities and Power for Exascale Computing are growing at an unsustainable rate

## Open Science Center (40k ft2)

- Upgrading building power to 15 MW
- 210 MW substation, upgradeable to 280 MW
- Deploying a 6,600 ton chiller plant
- Tripling UPS and generator capability

## National Security Center (32k ft2)

- Capability computing for national defence
- 25 MW of power and 8,000+ ton chillers

## New Computer Facility (260k ft2)

- 110 ft2 raised floor classified; same unclassified
- Shared mechanical and electrical
- Lights out facility
- Capable of greater than 100 MW power

# Power, cooling, and floorspace are fundamental challenges for the entire community over the next decade

ExaScale Computing Study:
Technology Challenges in
Achieving Exascale Systems

Peter Kogge, Editor & Study Lead
Keren Bergman
Shekhar Borkar
Dan Campbell
William Carlson
William Dally
Monty Denneau
Paul Franzon
William Harrod
Kerry Hill
Jon Hiller
Sherman Karp
Stephen Keckler
Dean Klein
Robert Lucas
Mark Richards
Al Scarpelli
Steven Scott
Allan Snavely
Thomas Sterling
R. Stanley Williams
Katherine Yelick

September 28, 2008

This work was sponsored by DARPA IPTO in the ExaScale Computing Study with Dr. William Harrod as Program Manager; AFRL contract number FA8650-07-C-7724. This report is published in the interest of scientific and technical information exchange and its publication does not constitute the Government's approval or disapproval of its ideas or findings

## NOTICE

Using Government drawings, specifications, or other data included in this document for any purpose other than Government procurement does not in any way obligate the U.S. Government. The fact that the Government formulated or supplied the drawings, specifications, or other data does not license the holder or any other person or corporation; or convey any rights or permission to manufacture, use, or sell any patented invention that may relate to them.

APPROVED FOR PUBLIC RELEASE, DISTRIBUTION UNLIMITED.

- **Energy, Power, and Facilities**
  - > Exascale system cannot use 1 GW!!
  - > Can we do Exascale in 20 MW??

- **Memory and Storage**

- **Concurrency and Locality**

- **Resiliency**

# Decisions for Managing Power and Facilities Demands

- **Build bigger buildings near power stations, and allocate large budget for power**



- **Improve efficiencies**
  - > Power distribution
  - > Workload scheduling
  - > Software
- **Design new underlying technologies**
  - > Optical networks
  - > 3D stacking
  - > MRAM, PCM, nanotubes
- **Use different architectures (that match your workload)**
- **Or, build bigger buildings and plan to pay $$$ for ops**

# Scaling dimension seems straightforward …

240 cores, 30k threads…



160…



80…



4…

# Another dimension:
# Several Factors Motivate Heterogeneity

- **Sacrifice generality and compatibility to address specific algorithms**
- **Computational power**
  - > Asymmetric, dynamic cores can provide performance advantages [1]
  - > Use transistors differently
- **Tukwila – First 2 billion transistor chip**
  - > 80486 had ~1.2M transistors, ~50MHz, 1989
  - > Specialization can be free

- **Power**
  - > Combination of these features provide more ops per watt for your targeted application
- **Memory**
  - > Include different cache or memory organization tuned for data access patterns
- **Performance of general purpose processors de-accelerating**



ISSCC 2008.

[1] M.D. Hill and M.R. Marty, "Amdahl's Law in the Multicore Era," *IEEE Computer*, to appear, 2008.

# *ENABLING HETEROGENEOUS COMPUTING FOR THE OPEN SCIENCE COMMUNITY*

# National Science Foundation Office of Cyberinfrastructure

- NSF OCI HPC Strategy
  - Track 1, BlueWaters
  - Track 2, TACC Ranger, NICS Kraken
- NSF 08-573 OCI Track 2D RFP in Fall 2008
  - Data Intensive
  - Experimental Grid testbed
  - Pool of loosely coupled grid-computing resources
  - Experimental HPC System of Innovative Design

**An experimental high-performance computing system of innovative design.** Proposals are sought for the development and deployment of a system with an architectural design that is outside the mainstream of what is routinely available from computer vendors. Such a project may be for a duration of up to five years and for a total award size of up to $12,000,000. It is not necessary that the system be deployed early in the project; for example, a lengthy development phase might be included. Proposals should explain why such a resource will expand the range of research projects that scientists and engineers can tackle and include some examples of science and engineering questions to which the system will be applied. It is not necessary that the design of the proposed system be useful for all classes of computational science and engineering problems. When finally deployed, the system should be integrated into the TeraGrid. It is is anticipated that the system, once deployed, will be an experimental TeraGrid resource, used by a smaller number of researchers than is typical for a large TeraGrid resource. (Up to 5 years duration. Up to $12,000,000 in total budget to include development and/or acquisition, operations and maintenance, including user support. First-year budget not to exceed $4,000,000.)

Georgia Tech | NICS | THE UNIVERSITY of TENNESSEE | OAK RIDGE National Laboratory | NVIDIA | hp | NSF

# Large Scale Computational Science Demands Performance, Programmability, Precision, Reliability, Cost from HPC Platforms

- Performance
  - Must show reasonable performance improvements at scale on real scientific applications of interest
- Programmability
  - Must be easy to re-port and re-optimize applications for each new architecture (generation) without large effort, delays
- Precision - Accuracy
  - Must provide impressive performance accurately
- Reliability
  - Must get high scientific throughput without job failures or inaccurate results
- Power and Facilities Cost
  - Must be reasonably affordable in terms of power and facilities costs

# Oct 2008 – Alternatives Analysis

- STI Cell
- FGPAs
- Cyclops64
- Cray XMT
- Sun Rock/Niagara
- ClearSpeed

- Tensilica
- Tilera
- Anton
- SGI Molecule
- Intel Larrabee
- Graphics processors

- Others…

# Alternatives analysis concluded GPUs were a competitive solution

- Success with various applications at DOE, NSF, government, industry
  - Signal processing, image processing, etc.
  - DCA++, S3D, NAMD, many others
- Commodity solution
  - Certified GPGPU clusters widely available this past quarter from multiple vendors
- Improving programmability
  - Widespread adoption of CUDA
  - OpenCL ratification
- Community application experiences also positive
  - Frequent workshops, tutorials, software development, university classes
  - Many apps teams are excited about using GPGPUs

- NVIDIA GT200 - 240

# GPU Rationale – What's different now?

# GPU Rationale - Power Efficiency!

# Keeneland – An NSF-Funded Partnership to Enable Large-scale Computational Science on Heterogeneous Architectures

- Track 2D System of Innovative Design
  - Large GPU cluster
    - Initial delivery system – Spring 2010
    - Full scale system – Spring 2012
- Software tools, application development
- Operations, user support
- Education, Outreach, Training for scientists, students, industry

| | |
|---|---|
| Jeffrey | Vetter |
| Jack | Dongarra |
| Richard | Fujimoto |
| Thomas | Schulthess |
| Karsten | Schwan |
| Sudha | Yalamanchili |
| Kathlyn | Boudwin |
| Jim | Ferguson |
| Doug | Hudson |
| Patricia | Kovatch |
| Bruce | Loftis |
| Jeremy | Meredith |
| Jim | Rogers |
| Philip | Roth |
| Arlene | Washington |
| Phil | Andrews |
| Mark | Fahey |
| Don | Reed |
| Tracy | Rafferty |
| Ursula | Henderson |
| Terry | Moore |

# Keeneland Partners

| Georgia Institute of Technology | UT National Institute of Computational Sciences | Oak Ridge National Laboratory | University of Tennessee, Knoxville | NVIDIA | HP |
|---|---|---|---|---|---|
| Project management | Operations and TG/XD Integration | Applications | Scientific Libraries | Tesla | HPC Host System |
| Acquisition and alternatives assessment | User and Application Support | Facilities | Education, Outreach, Training | Applications optimizations | System integration |
| System software and development tools | Operational Infrastructure | Education, Outreach, Training | | Training | Training |
| Education, Outreach, Training | Education, Outreach, Training | | | | |

# Keeneland Initial Delivery (ID) System

- Hewlett Packard Nodes
  - Dual socket Intel 2.8 GHz Nehalem-EP
  - 24 GB Main memory per node
- NVIDIA Servers
  - Fermi GPUs
- InfiniBand 4x QDR w/ full bisection interconnect
- Traditional Linux software stack augmented with GPU compilers, software tools, libraries
- Hundreds of Fermi processors
- Delivery and acceptance in Spring 2010

# ID system will use NVIDIA's Fermi

- "The soul of a supercomputer in the body of a GPU."
- 3B transistors
- ECC
- 8x the peak double precision arithmetic performance over NVIDIA's last generation GPU.
- 512 CUDA Cores featuring the new IEEE 754-2008 floating-point standard
- NVIDIA Parallel DataCache
- NVIDIA GigaThread Engine
- CUDA and OpenCL support
- Debuggers, language support

# *APPLICATIONS*

# Computational Materials - Case Study

- Quantum Monte Carlo simulation
  - High-temperature superconductivity and other materials science
  - 2008 Gordon Bell Prize
- GPU acceleration speedup of 19x in main QMC Update routine
  - Single precision for CPU and GPU: target single-precision only cards
  - Required detailed accuracy study and mixed precision port of app
- Full parallel app is 5x faster, start to finish, on a GPU-enabled cluster

GPU study: J.S. Meredith, G. Alvarez, T.A. Maier, T.C. Schulthess, J.S. Vetter, "Accuracy and Performance of Graphics Processors: A Quantum Monte Carlo Application Case Study", *Parallel Comput., 35(3):151-63, 2009.*

Accuracy study: G. Alvarez, M.S. Summers, D.E. Maxwell, M. Eisenbach, J.S. Meredith, J. M. Larkin, J. Levesque, T. A. Maier, P.R.C. Kent, E.F. D'Azevedo, T.C. Schulthess, "New algorithm to enable 400+ TFlop/s sustained performance in simulations of disorder effects in high-Tc superconductors", SuperComputing, 2008. [*Gordon Bell Prize winner*]

Chart: CPU Runtime (red) and GPU Runtime (green) vs. number (4, 8, 12, 16, 20, 24); y-axis from 1 seconds to 10000 seconds.

# Combustion with S3D – Case Study

- Application for combustion - S3D
  - Massively parallel direct numerical solver (DNS) for the full compressible Navier-Stokes, total energy, species and mass continuity equations
  - Coupled with detailed chemistry
  - Scales to 150k cores on Jaguar

- Accelerated version of S3D's Getrates kernel in CUDA
  - 14.3x SP speedup
  - 9.32x DP speedup



K. Spafford, J. Meredith, J. S. Vetter, J. Chen, R. Grout, and R. Sankaran. Accelerating S3D: A GPGPU Case Study. Proceedings of the Seventh International Workshop on Algorithms, Models, and Tools for Parallel Computing on Heterogeneous Platforms (HeteroPar 2009) Delft, The Netherlands.

# Biomolecular systems from NAMD Team – Not just us

- NAMD, VMD
  - Study of the structure and function of biological molecules

- Calculation of non-bonded forces on GPUs leads to 9x speedup

- Framework hides most of the GPU complexity from users



J.C. Phillips and J.E. Stone, "Probing biomolecular machines with graphics processors," *Commun. ACM*, 52(10):34-41, 2009. (fig)

Georgia Tech    NICS    THE UNIVERSITY of TENNESSEE    OAK RIDGE National Laboratory    NVIDIA    hp    NSF

# KEENELAND SOFTWARE

# Keeneland Software Environment

- Integrated with NSF TeraGrid/XD
  - Including TG and NICS software stack

- Programming models
  - CUDA
  - OpenCL
  - PGI w/ accelerate
  - OpenMP 3.0
  - MPI

- Additional software activities
  - Performance and correctness tools
  - Scientific libraries
  - Virtualization
  - Benchmarks

Georgia Tech    NICS    THE UNIVERSITY of TENNESSEE    OAK RIDGE National Laboratory    nVIDIA    hp    NSF
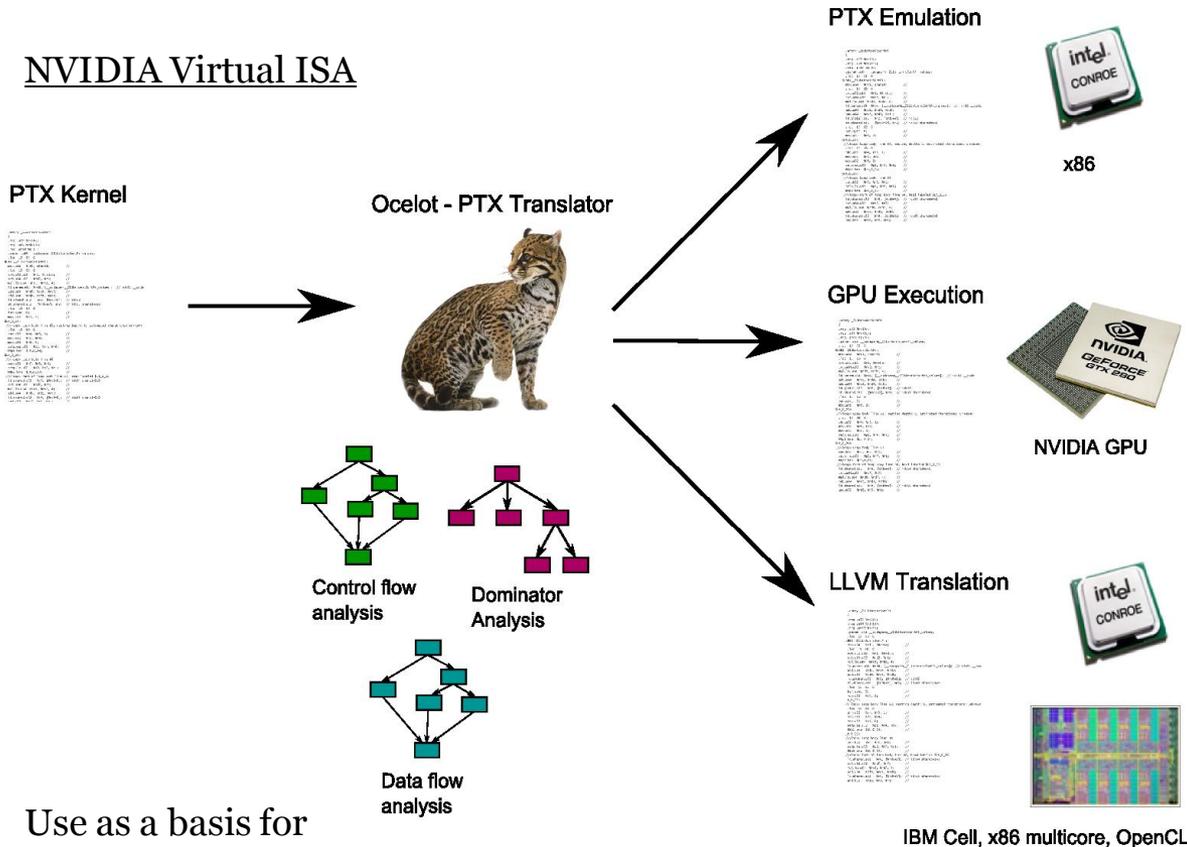
# Ocelot: Dynamic Execution Infrastructure

NVIDIA Virtual ISA

PTX Kernel

Ocelot - PTX Translator

PTX Emulation

x86

PTX 1.4 compliant Emulation
• Validated on full CUDA SDK
• Open Source version released

*http://code.google.com/p/gpuocelot/*

Control flow analysis

Dominator Analysis

Data flow analysis

GPU Execution

NVIDIA GPU

LLVM Translation

IBM Cell, x86 multicore, OpenCL

*Gregory Diamos, Dhuv Choudhary, Andrew Kerr, Sudhakar Yalamanchili*

*Gregory Diamos recently awarded NVIDIA Fellowship*

Use as a basis for
- Insight → workload characterization
- Performance tuning → detecting memory bank conflicts
- Debugging → illegal memory accesses, out of bounds checks, etc.

Georgia Tech   NICS   THE UNIVERSITY of TENNESSEE   OAK RIDGE National Laboratory   NVIDIA   hp   NSF

# Workload Analysis: Examples



**Branch Divergence**
• Study of control Flow behavior
• Motivate synchronization support

**Inter-thread Data Flow**
• Study of data sharing patterns
• Motivate architectural support

*Gregory Diamos, Dhuv Choudhary, Andrew Kerr, Sudhakar Yalamanchili*

# One and two-sided Multicore+GPU Factorizations

- These will be included in up-coming MAGMA releases

- Two-sided factorizations can not be efficiently accelerated on homogeneous x86-based multicores (above) because of memory-bound operations
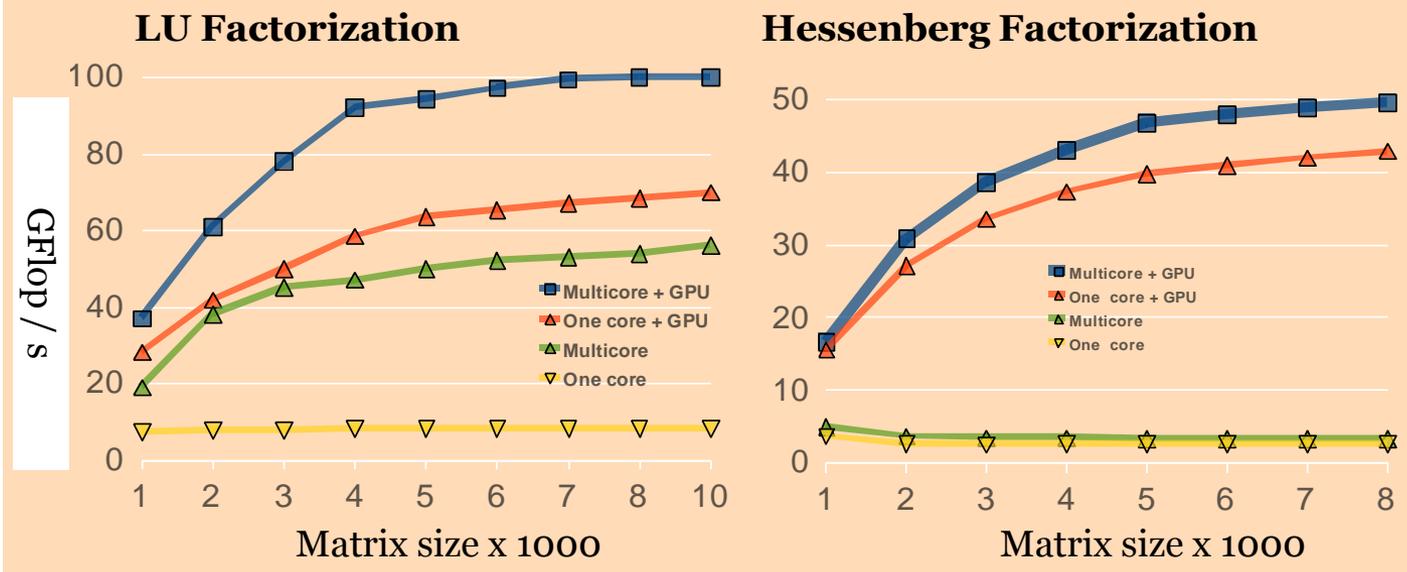  - MAGMA provided hybrid algorithms that overcome those bottlenecks (16x speedup!)



**Multicore + GPU Performance in double precision**

*LU Factorization* — Multicore + GPU, One core + GPU, Multicore, One core; GFlop / s vs Matrix size x 1000

*Hessenberg Factorization* — Multicore + GPU, One core + GPU, Multicore, One core; Matrix size x 1000

*Jack Dongarra, Stan Tomov, and Rajib Nath*
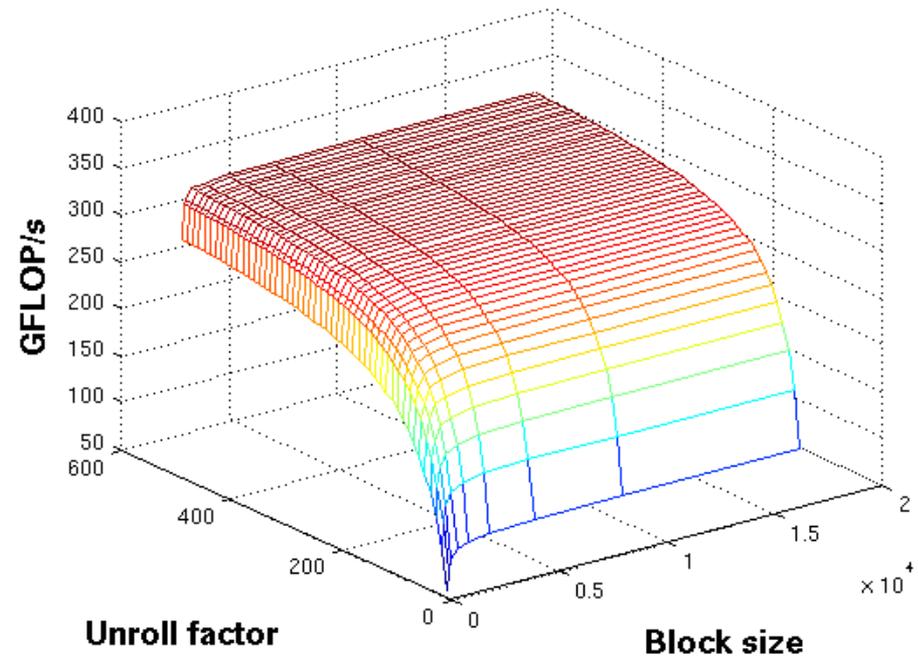
**GPU** : NVIDIA GeForce GTX 280
**CPU** : Intel Xeon dual socket quad-core @2.33 GHz
**GPU BLAS** : CUBLAS 2.2, dgemm peak: 75 GFlop/s
**CPU BLAS** : MKL 10.0 , dgemm peak: 65 GFlop/s

# Scalable Heterogeneous Computing (SHOC) Benchmark Suite

- Problem: we don't yet understand overheads and performance benefits of GPU versus CPU, and trade-offs between different programming approaches
    - Systems can differ greatly
    - Need procurement and acceptance tests

- Antero: a collection of benchmark programs for quantifying those overheads and trade-offs

- Multiple levels
    - Level 0: low-level, feeds and speeds
    - Level 1: higher-level, problem kernels

- Multiple test categories
    - Performance
    - Stability

- Support for both single-node and clusters like Keeneland (MPI-based)



Jeremy Meredith, Phil Roth, Kyle Spafford, Anthony Danalis, Gabriel Marin, Vinod Tipparaju, Jeffrey Vetter

Georgia Tech    NICS    THE UNIVERSITY of TENNESSEE    OAK RIDGE National Laboratory    NVIDIA    hp    NSF

# Summary

- Predictive scientific simulation is important for scientific discovery
  - Advancing science, informing policymakers
- HPC systems have been highly successful
  - Decades of improvement
- The HPC community has several (new) constraints
  - Power, Facilities, Cost
- Heterogeneous computing with GPUs offers some opportunities and challenges
  - High performance; good performance per watt
  - Programmability; limited applicability
- Keeneland - Newly awarded NSF partnership will provide heterogeneous supercomputing for open science

# *Thank You!*

**Thanks to contributors, sponsors**

Many collaborators across apps teams, academia, labs, and industry

DOE, NSF, ORNL, DARPA, DOD

**More information**

http://ft.ornl.gov

vetter@computer.org

Publications: http://ft.ornl.gov/pubs

**http://keeneland.gatech.edu**

**http://www.cse.gatech.edu**

**http://www.cercs.gatech.edu**

**http://icl.cs.utk.edu**

**http://www.nics.tennessee.edu/**

**http://ft.ornl.gov**

**http://nsf.gov/dir/index.jsp?org=OCI**