# Data Analysis and Visualization at Exascale

**Dr. Lucy Nowell**

**Computer Scientist and Program Manager**

**Office of Advanced Scientific Computing Research**

**DOE Office of Science**

**SC 10 – November 2010**

# Leadership Computing Facilities
### *The Office of Science leads the World in supercomputing capabilities*



"Supercomputer modeling and simulation are changing the face of science and sharpening America's competitive edge."

Secretary Steven Chu

The Cray XT5 Supercomputer at Oak Ridge National Lab can perform over 2.3 quadrillion operations per second. It ranks #1 of the fastest computers world wide by Top500.org

# ASCR Computer Science
# Base Research

- **ASCR Base CS Program tries to address two fundamental questions:**
  - How can we make today's and tomorrow's leading edge computers tools for science?
  - How do we extract scientific information from extreme scale data from experiments and simulation?

- **There are several factors that provide important context for the ASCR Base CS program:**
  - SciDAC Centers and Institutes
  - Research and Evaluation Partnerships
  - ASCR Facilities

# Increasing Machine Capability

- **Gigaflop = one billion (1,000,000,000,000) floating point operations (flops) per second**
- **Teraflop = ~1024 gigaflops, or roughly 1 trillion flops**
- **Petaflop = ~1 quadrillion (or $10^{15}$) flops, or 1024 teraflops**
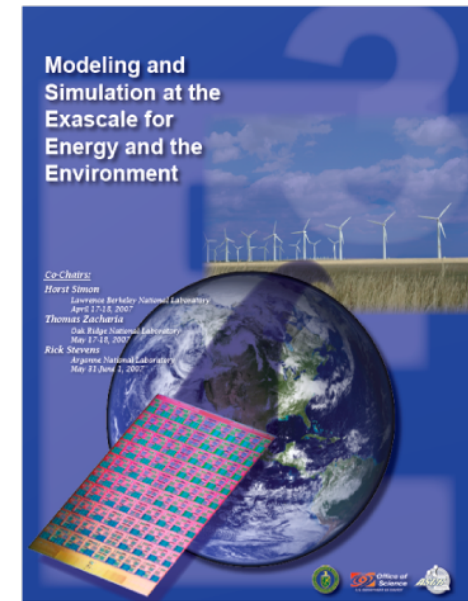- **Exaflop = 1 quintillion (or $10^{18}$) flops, or 1 million teraflops**

# What Was That Again?

- **7 Gigaflops = O(1 floating point operations per second for every person on Earth)**

- **7 Teraflops = O(~1,000 flops for every person on Earth)**

- **7 Petaflops = O(1 million flops for every person on Earth)**

- **7 Exaflops = O(1 billion flops for every person on Earth)**
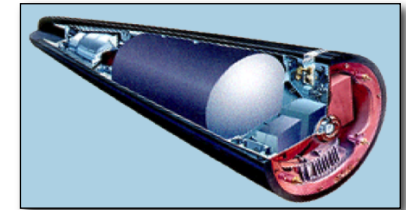
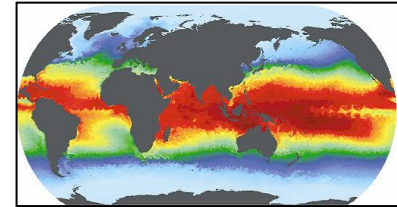# Why Do We Need a Machine That Large?

- **It will take a machine capable of ~3 Exaflops to completely model the U.S. power grid. – Steve Elbert, PNNL**

- **An Exascale computer will be able to model weather/climate at a resolution of roughly 10 meters.**

**U.S. DEPARTMENT OF ENERGY**
Office of Science

# DOE mission imperatives require simulation and analysis for policy and decision making
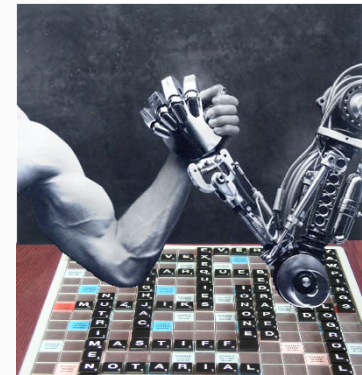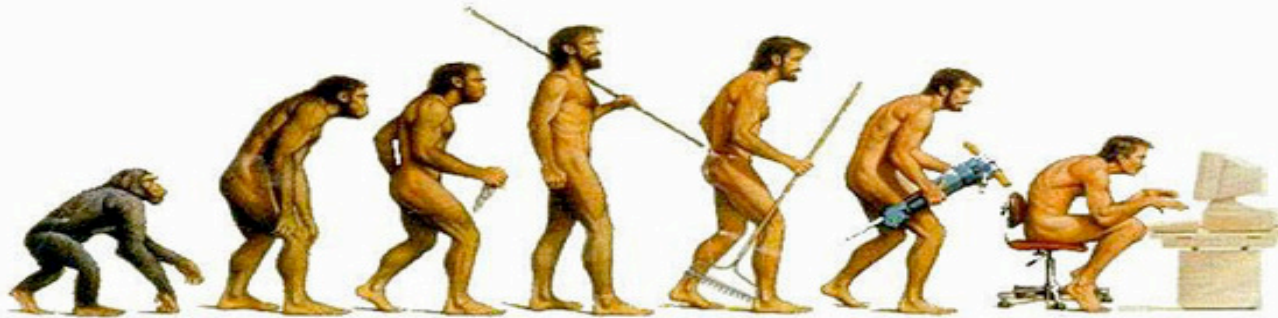
- *Climate Change*: **Understanding and mitigating the effects of global warming**
  - **Sea level rise**
  - **Severe weather**
  - **Regional climate change**
  - **Geologic carbon sequestration**
- *National Nuclear Security*: **Maintaining a safe, secure and reliable nuclear stockpile**
  - **Stockpile certification**
  - **Predictive scientific challenges**
  - **Real-time evaluation of urban nuclear detonation**
- *Energy*: **Reducing U.S. reliance on foreign energy sources and reducing the carbon footprint of energy production**
  - **Reducing time and cost of reactor design and deployment**
  - **Improving the efficiency of combustion energy sources**

Accomplishing these missions requires exascale resources.

# Evolution to Exascale



Gigascale          Terascale          Petascale          Exascale

# (Exa)Scale Changes Everything

|  | 2010 | 2018 | Factor Change |
|---|---|---|---|
| System peak | 2 Pf/s | 1 Ef/s | 500 |
| Power | 6 MW | 20 MW | 3 |
| System Memory | 0.3 PB | 10 PB | 33 |
| Node Performance | 0.125 Gf/s | 10 Tf/s | 80 |
| Node Memory BW | 25 GB/s | 400 GB/s | 16 |
| Node Concurrency | 12 cpus | 1,000 cpus | 83 |
| Interconnect BW | 1.5 GB/s | 50 GB/s | 33 |
| System Size (nodes) | 20 K nodes | 1 M nodes | 50 |
| Total Concurrency | 225 K | 1 B | 4,444 |
| Storage | 15 PB | 300 PB | 20 |
| Input/Output bandwidth | 0.2 TB/s | 20 TB/s | 100 |

DOE Exascale Initiative Roadmap, Architecture and Technology Workshop, San Diego, December, 2009.

# Exascale Challenges

## Exascale ≠ Petascale X 1000

- Total concurrency in the applications must rise by a factor of ~1 million;
- Memory per processor falls dramatically which makes current weak scaling approaches problematic;
- For both power and performance reasons, locality of data and computation is much more important
- The failure rates for components and manufacturing variability make it unreasonable to assume the computer is deterministic. This is true for performance today and will affect the results of computations by 2018 due to silent errors.
- Synchronization will be very expensive. In addition, work required to manage synchronization is high.
- The I/O system at all levels – chip to memory, memory to I/O node, I/O node to disk—  will be much harder to manage due to the relative speeds of the components.

# Science at Scale

- "From a scientist's perspective, the ratio of memory to processor is critical in determining the size of the problem that can be solved. Remember that the processor dictates how much computing can be done; the memory dictates the size of the problem that can be handled. In the Exascale design…there is 500 times more compute power, however only 30 times the memory, so applications cannot just scale to the speed of the machine. Scientists and computer scientists will have to rethink how they are going to use these systems. This factor of >10 loss in memory/compute power means potentially totally redesigning the current application codes."

P.49 ASCAC Exascale report, October 2010

# Challenges for the Future
## Mountains of data

- **Storing:**
  - **Long term: where do we put 500TB?**
  - **Short term: scratch ~ 1TB, but need ~ 10TB!**
- **Moving:**
  - **Archive to scratch (~ 2 weeks to move 10TB)**
  - **HPC facility to local analysis cluster (longer)**
- **Processing:**
  - **Everything must be parallel, scalable.**
  - **IO speed, memory are the bottlenecks.**
- **Transforming Data into Insight**
  - **Physics are more complex**
  - **Wider range of scales, manual sifting is impossible.**
  - **Multi-scale analysis methods**
  - **Feature detection, growing, and tracking**

HPSS storage facility at NERSC



*"Where is the wisdom that is lost in knowledge? Where is the knowledge we have lost in information?"*

*-T.S. Eliot*

# Data Integration Imperative

- **"We seek solutions. We don't seek – dare I say this? – just scientific papers anymore."**
  - **Secretary Chu,** http://articles.sfgate.com/2007-03-22/bay-area/17236680_1_uc-berkeley-alternative-fuels-fossil-fuels

- **Science challenges require integration of data from multiple models**

  - **Climate modeling and the impact of climate change: ocean models, atmospheric and cloud models, ice sheet models, etc.**

  - **Implications of energy production: Will use of biomass lead to starving populations because of changes in crops or changed use of crops? Will it motivate increased burn-off of forests, with implications for climate change?**

13

# Heterogeneous Data: Environmental Science Example

- **Mobile stations**
- **High-resolution weather stations**
- **Full-size snow/weather stations**
- **External weather stations**
- **Satellite imagery**
- **Weather radar**
- **Mobile weather radar**
- **Stream observations**
- **Citizen-supplied observations**
- **Ground LIDAR**
- **Aerial LIDAR**

- **Nitrogen/methane measures**
- **Snow hydrology & avalanche probes**
- **Seismic probes**
- **Distributed optical fiber temperature sensing**
- **Water quality sampling**
- **Stream gauging stations**
- **Rapid mass movements research**
- **Run-off stations**
- **Soil research**

Source:  Lehning, Michael et al, "Instrumenting the Earth: Next-Generation Sensor Networks and Environmental Science" in The Fourth Paradigm: Data-Intensive Science, ed. Tony Hey, http://research.microsoft.com/en-us/collaboration/fourthparadigm/4th_paradigm_book_complete_lr.pdf

14

# Data Surpass
# Human Cognitive Limits

- **Most of us read at most a few gigabytes in a lifetime.**

- **We are not good at detecting patterns in large volumes of data.**

- **The data of science are mostly numeric, and humans process numbers even less well than text.**

- **Cognitive biases distort or inhibit understanding.**

# Visual Analysis

- **Transforms the cognitive problem to a perceptual one, taking advantage of the much broader bandwidth of human vision**

- **Should support human intuition and insight**
  - **Exploration of data and detection of patterns**
  - **Testing of hypotheses about the data**
  - **Requires deep integration with both data and the underlying analytic engine**

- **Requires deep knowledge of human perceptual and cognitive characteristics, as well as knowledge of the science application area**

- **Can make science understandable (and exciting) to policy makers and the general public**

**U.S. DEPARTMENT OF**
**ENERGY**

Office of Science

Thank you!

Lucy.Nowell@Science.DOE.gov