

# Leadership Computing Directions at Oak Ridge National Laboratory: Navigating the Transition to Heterogeneous Architectures

International Symposium on  
Computing in Atmospheric Sciences  
September 8-12 , 2013

***James J. Hack, Director***  
**National Center for Computational Sciences**  
**Oak Ridge National Laboratory**

*Arthur S. "Buddy" Bland, OLCF Project Director*

*Bronson Messer, OLCF Scientific Computing*

*James Rogers, NCCS Director of Operations*

*Jack Wells, NCCS Director of Science*



U.S. DEPARTMENT OF  
**ENERGY**

 **OAK RIDGE NATIONAL LABORATORY**

MANAGED BY UT-BATTELLE FOR THE DEPARTMENT OF ENERGY

# U.S. Department of Energy strategic priorities

## Incremental changes to existing technologies cannot meet these challenges

- Transformational advances in energy technologies are needed
- Transformational adaptation strategies will need to be implemented
- Transformational changes to tools that allow a predictive explorations of paths forward

### Innovation

Investing in science, discovery and innovation to provide solutions to pressing energy challenges

### Energy

Providing clean, secure energy and promoting economic prosperity through energy efficiency and domestic forms of energy

### Security

Safeguarding nuclear and radiological materials, advancing responsible legacy cleanup, and maintaining nuclear deterrence





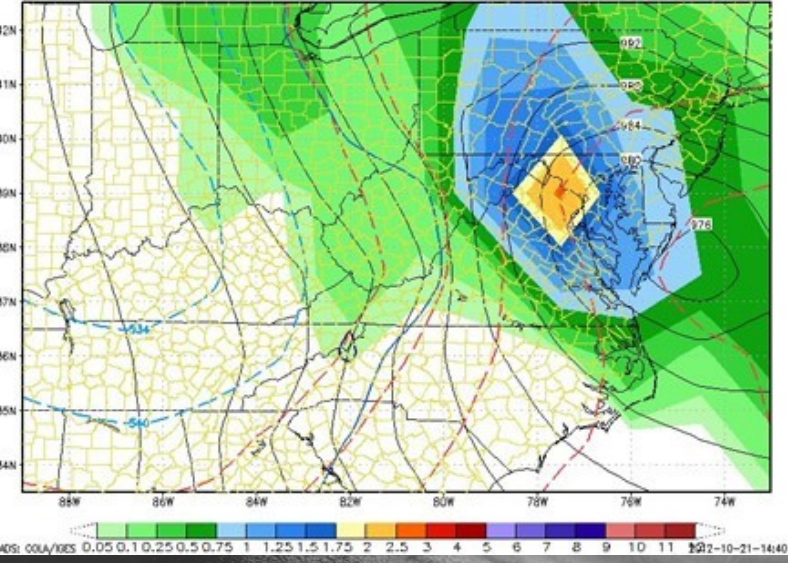
***New York Magazine***



# Example of the virtualization challenge

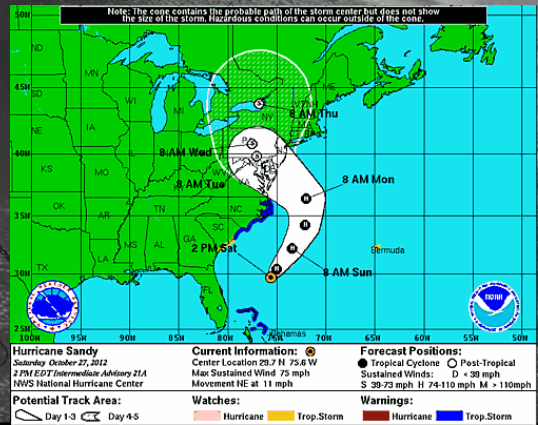
## And a great example of success

ECMWF MSLP, 1000-500 THK, 6 HR QPF Stormvistawxmodels.com  
Hour 186 (Mon 06Z29OCT2012)



*Big storm early next week (Oct 29) with wind and rain??? A number of computer models today were hinting at that, suggesting a tropical system now in the Caribbean may merge with a cold front in that time table, close to the East Coast. At the same time, 8 days is an eternity in storm forecasting, so plenty of time to watch this. Consider this an early “heads-up” that we’ll have an interesting an weather feature to monitor this week... details on what to expect and whether we’ll even deal with a storm still way up in the air.*

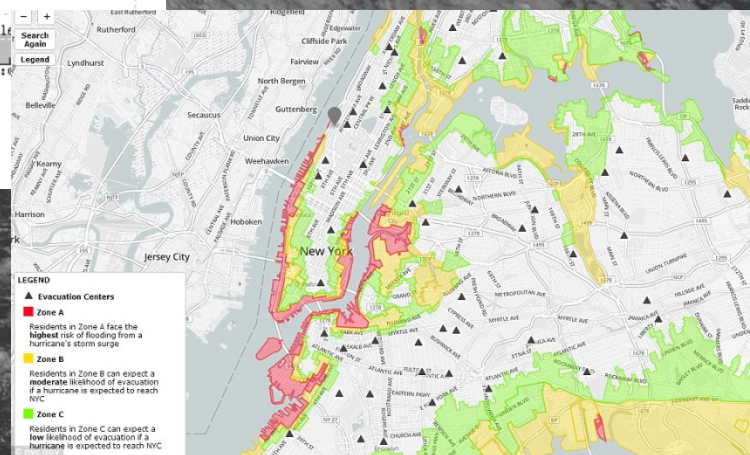
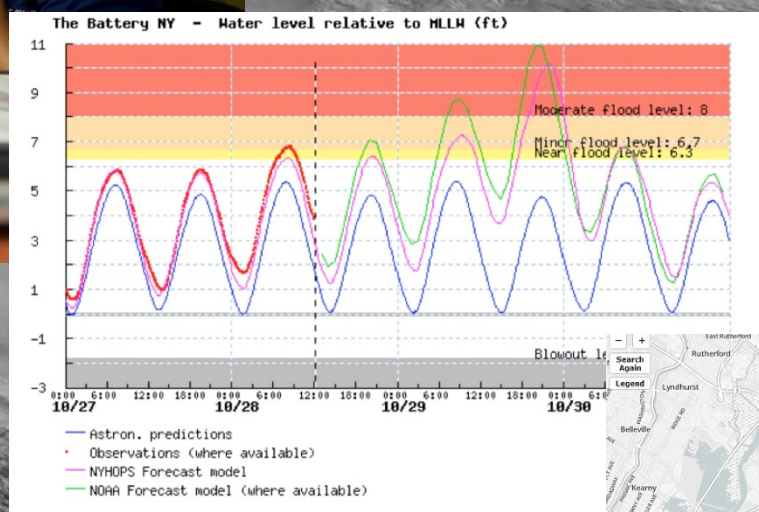
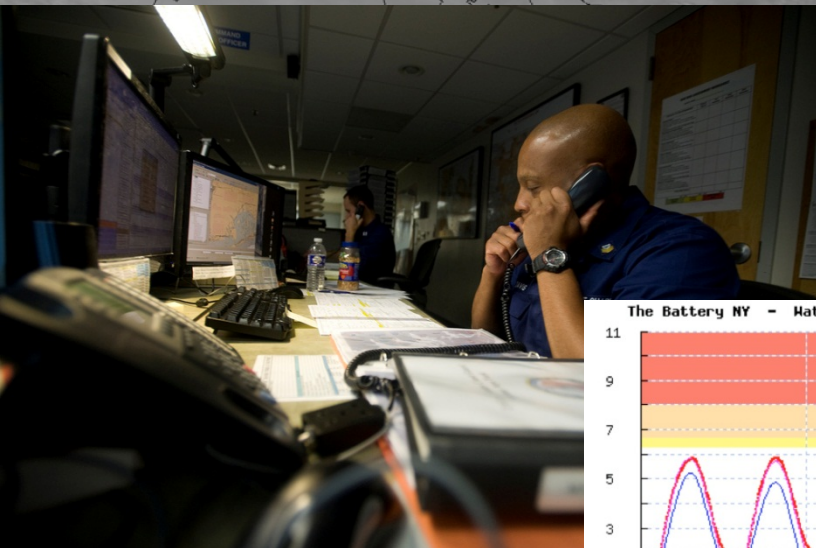
Jason Samenow, Capital Weather Gang, Washington Post 21 October 2012





# Example of the virtualization challenge

Simulation science has begun to master deterministic time scales and regional space scales



Also a demonstration of the scale challenge



# Example of the virtualization challenge

***Surprise!??***

Sweet Tweets

**CAN'T BELIEVE  
METEOROLOGISTS USED  
MATH AND SCIENCE TO  
PREDICT THIS STORM.  
THEY MUST BE MAGIC  
WIZARDS.**

***-Nate Silver***

[twitter.com/fivethirtyeight](https://twitter.com/fivethirtyeight)



***More of a Surprise!??***

***There has been a series of extreme weather incidents. That is not a political statement, that is a factual statement," Cuomo said. "Anyone who says there is not a dramatic change in weather patterns is denying reality."***

***"I said to the president kiddingly the other day, 'We have a one-hundred year flood every two years.'"***

***Governor Andrew Cuomo, 30 October 2012***

# Examples of climate consequences questions

- **Water Resources**

- management and maintenance of existing water supply systems, development of flood control systems and drought plans

- **Agriculture and food security**

- Erosion control, dam construction (irrigation), optimizing planting/harvesting times, introduction of tolerant/resistant crops (to drought, insect/pests, etc.)

- **Human health**

- Public health management reform, improved urban and housing design, improved disease/vector surveillance and monitoring

- **Terrestrial ecosystems**

- Improvement of management systems (deforestation, reforestation,...), development/improvement of forest fire management plans

- **Coastal zones and marine ecosystems**

- Better integrated coastal zone planning and management

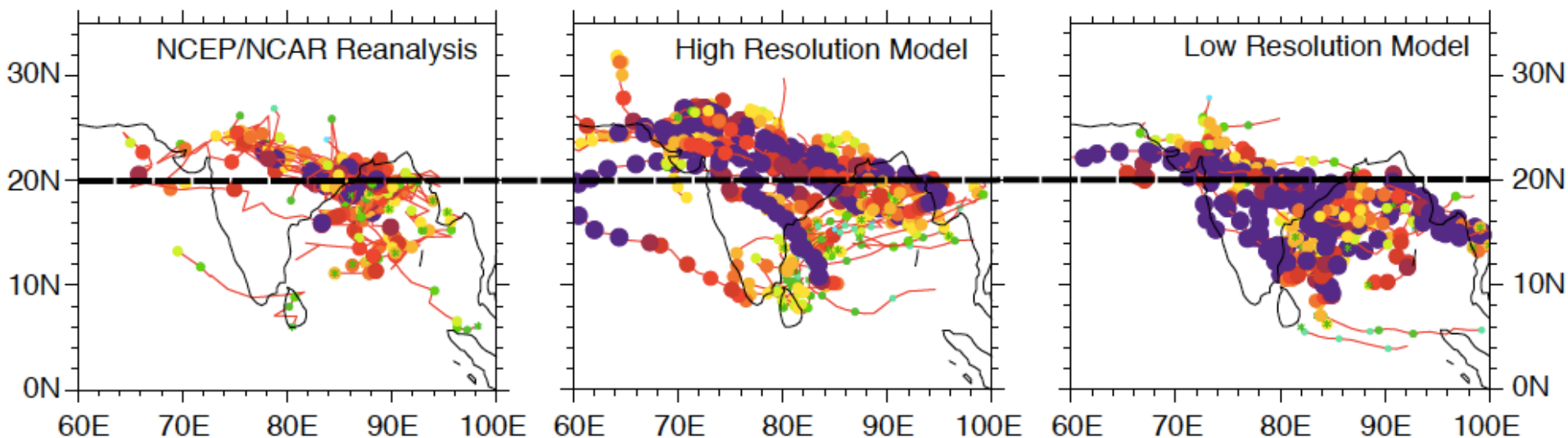
- **Human-engineered systems**

- Better planning for long-lived infrastructure investments

# South Asian Monsoon Depressions

*Ashfaq and collaborators*

## South Asian Summer Monsoon Depressions in 1999-2009



**Movement of monsoon depressions over land improves in the high-resolution model**



# High Resolution Spectral Configuration of CAM4

**Mahajan and collaborators**

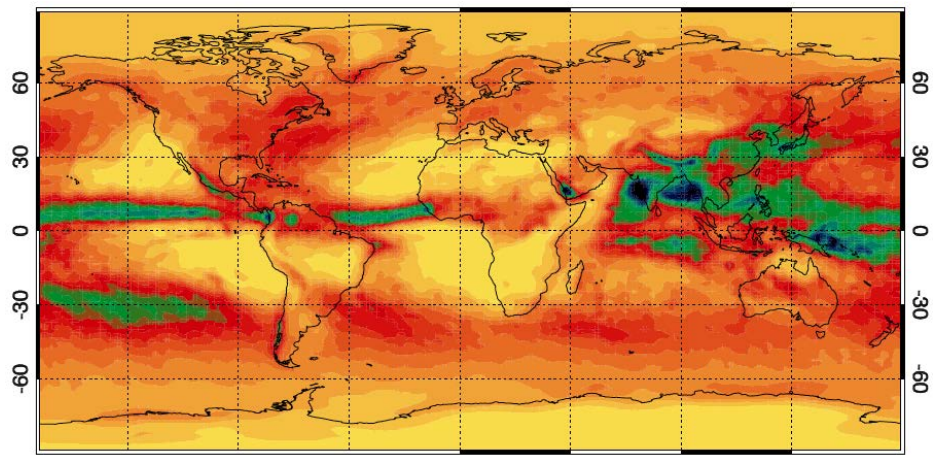
Experiments with T341 Spectral Configuration of CAM4 (ORNL):

- AMIP (1979-2005)
- CAM4 stand-alone pre-industrial control
- Fully Coupled pre-industrial control
- Fully Coupled present day

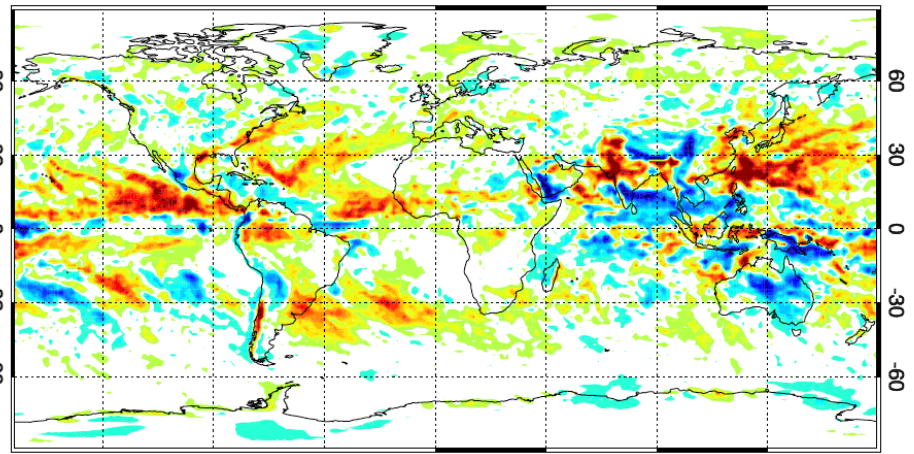
**Preliminary Results:**

\* General increase in variance of precipitation

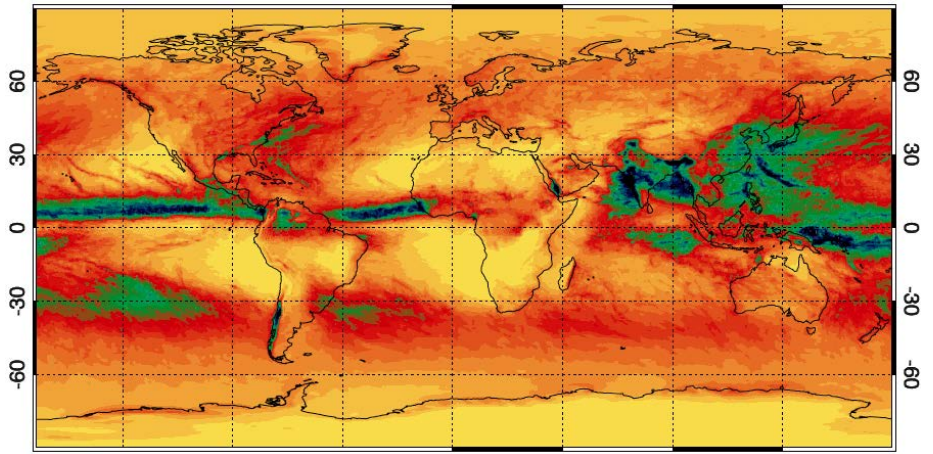
T85 Model (AMIP) Standard Deviation: Precipitation



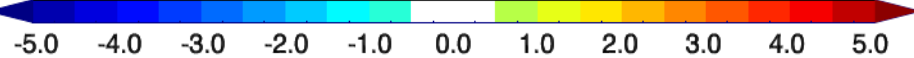
T341 Model – T85 Model



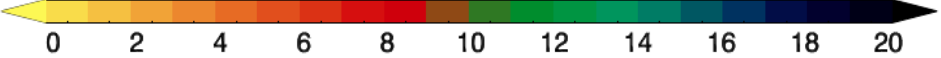
T341 Model (AMIP) Standard Deviation: Precipitation



Diff. in Standard Deviation: Precipitation (mm/day)



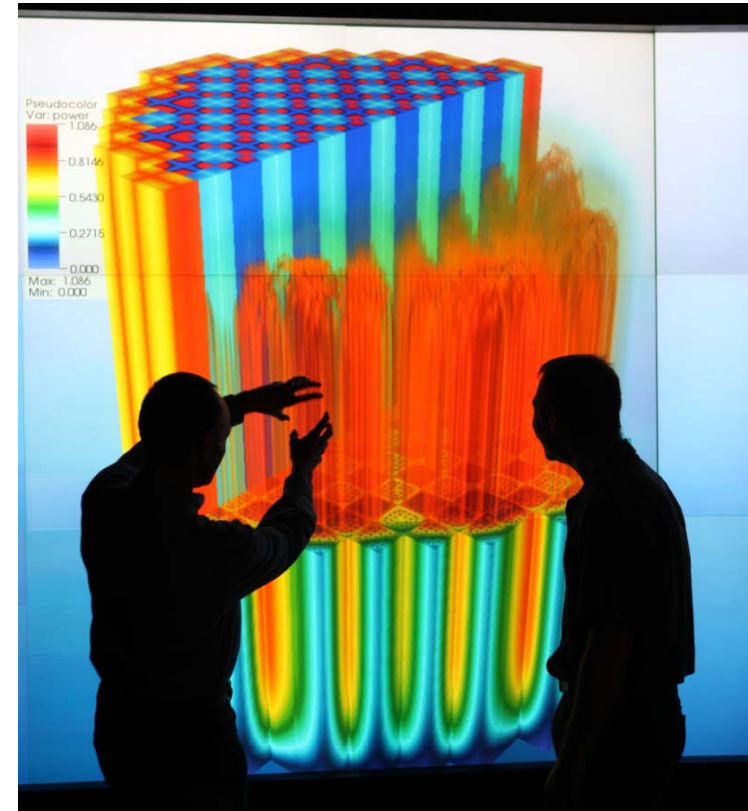
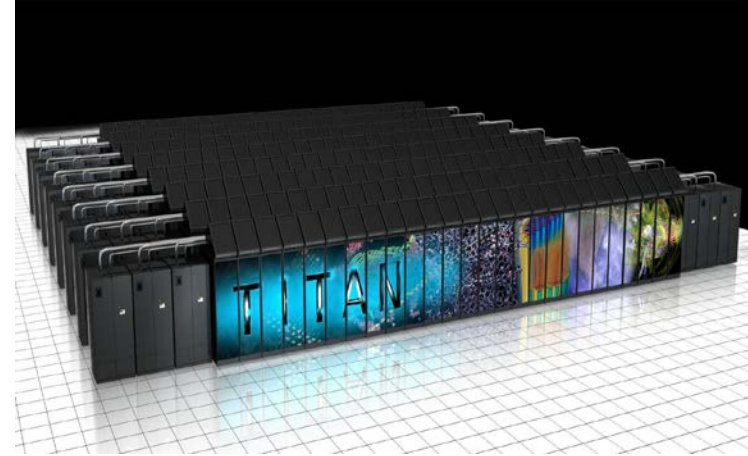
Standard Deviation: Precipitation (mm/day)



# Oak Ridge Leadership Computing Facility Mission

The OLCF is a DOE Office of Science National User Facility whose mission is to enable breakthrough science by:

- Fielding the most powerful capability computers for scientific research,
- Building the required infrastructure to facilitate user access to these computers,
- Selecting a few time-sensitive problems of national importance that can take advantage of these systems,
- And partnering with these teams to deliver breakthrough science.





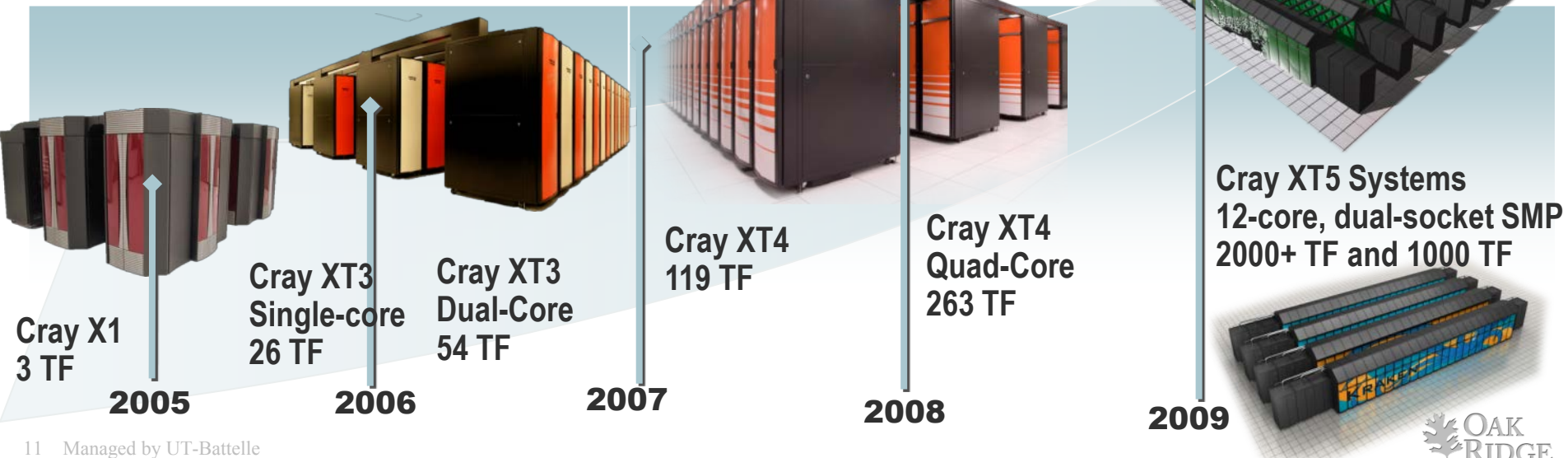
# ORNL increased system performance by 1,000 times between 2004 and 2009

Hardware scaled from single-core through dual-core to quad-core and dual-socket, 12-core SMP nodes

- NNSA and DoD have funded much of the basic system architecture research
  - Cray XT based on Sandia Red Storm
  - IBM BG designed with Livermore
  - Cray X1 designed in collaboration with DoD

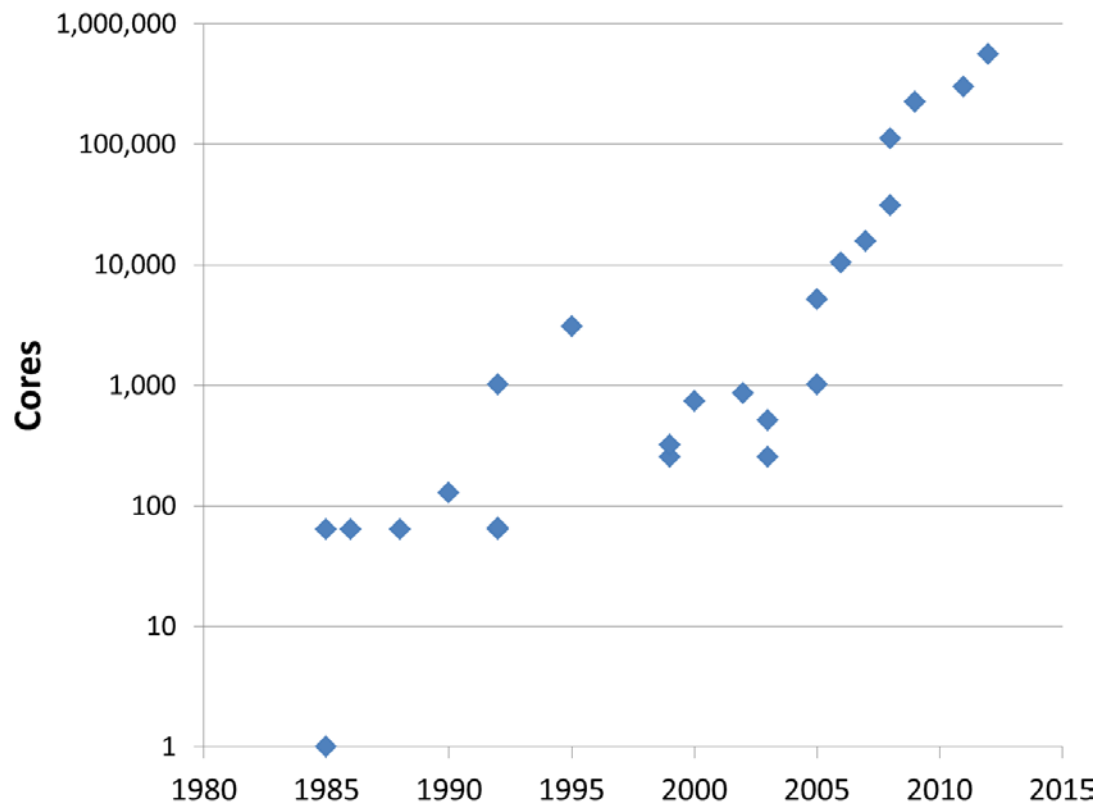
Scaling applications and system software is the biggest challenge

- DOE SciDAC and NSF PetaApps programs are funding scalable application work, advancing many apps
- DOE-SC and NSF have funded much of the library and applied math as well as tools
- Computational Liaisons key to using deployed systems



# Scaling of ORNL's Systems 1985 - 2013

- In the last 28 years ORNL systems have scaled from 64 cores to hundreds of thousands of cores and millions of simultaneous threads of execution
  - Multiple hierarchical levels of parallelism
  - Hybrid processors and systems
- Almost 30 years of application development focused on finding ways to exploit that parallelism

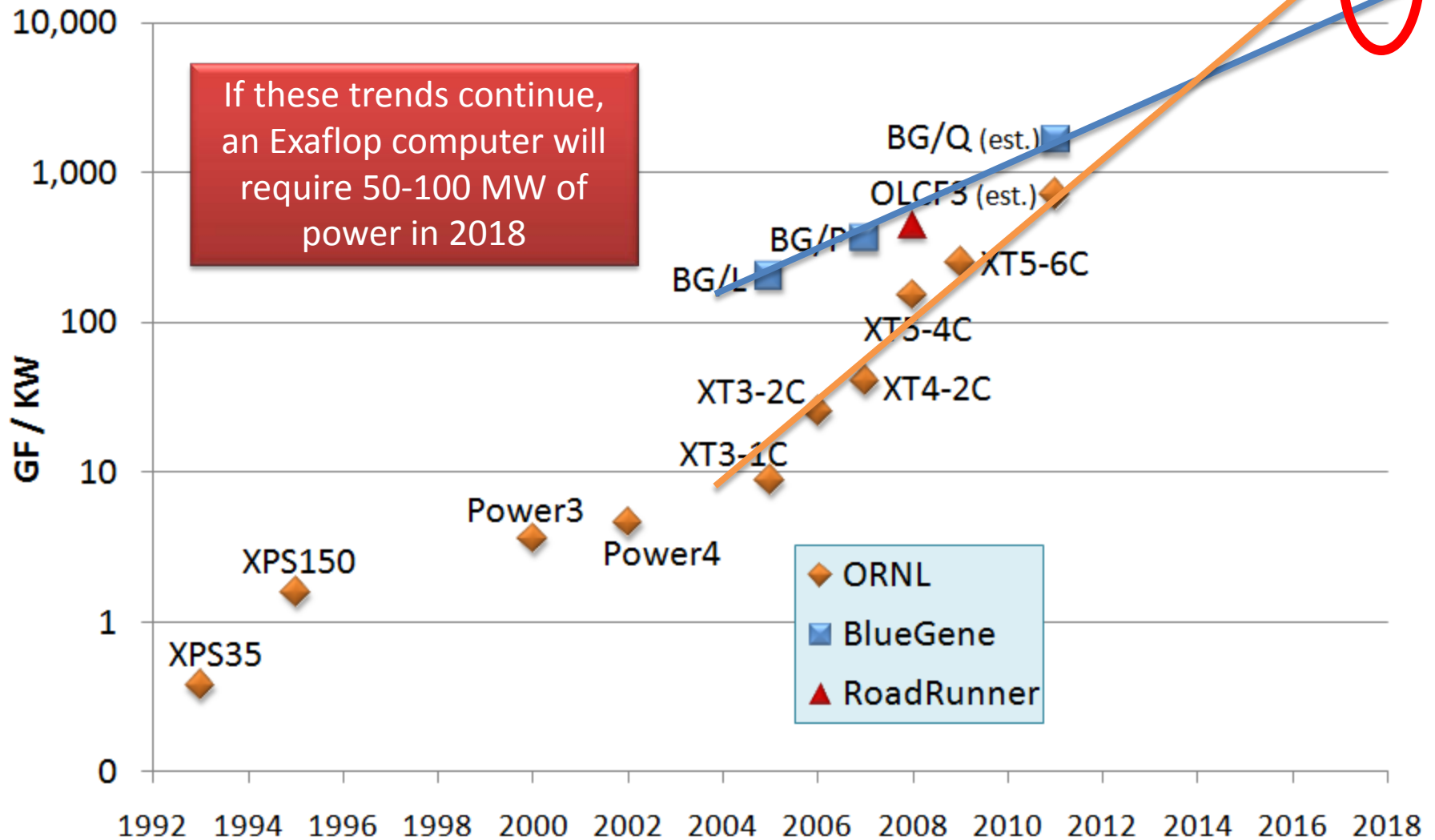




# Extreme-scale » Exascale Systems

- 1-10 billion way parallelism
  - Requires hierarchical parallelism to manage
    - MPI between nodes
    - OpenMP or other threads model on nodes
    - SIMD / Vector parallelism within thread
- Power will dominate architectures
  - Takes more power to go to memory for data than to recompute it
- Traditional “balance ratios” are eroding
  - Memory size is not growing as fast as processor performance
  - Memory bandwidth is growing even more slowly
  - Floating point operations cheap; memory access and data movement rate limitors

# Current technologies will require huge amounts of power for next generation systems



If these trends continue, an Exaflop computer will require 50-100 MW of power in 2018

◆ ORNL  
■ BlueGene  
▲ RoadRunner



# Technology transitions have been observed over time

Logistic change is characterized by an initial period of slow growth, followed by a period of exponential growth, then a point of inflection, and finally a period of asymptotic growth as the technology approaches a limit. This pattern of change was first observed in population studies [28], and it has since been found to be descriptive of change in a remarkably diverse set of circumstances, including technological evolution in general and the evolution of electronic and computer technologies in particular.

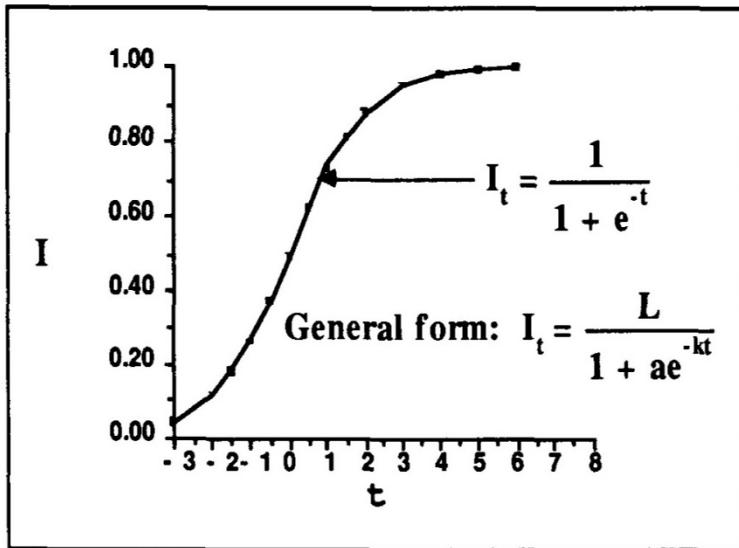


Figure 9. Logistic change.

Worlton (1988)

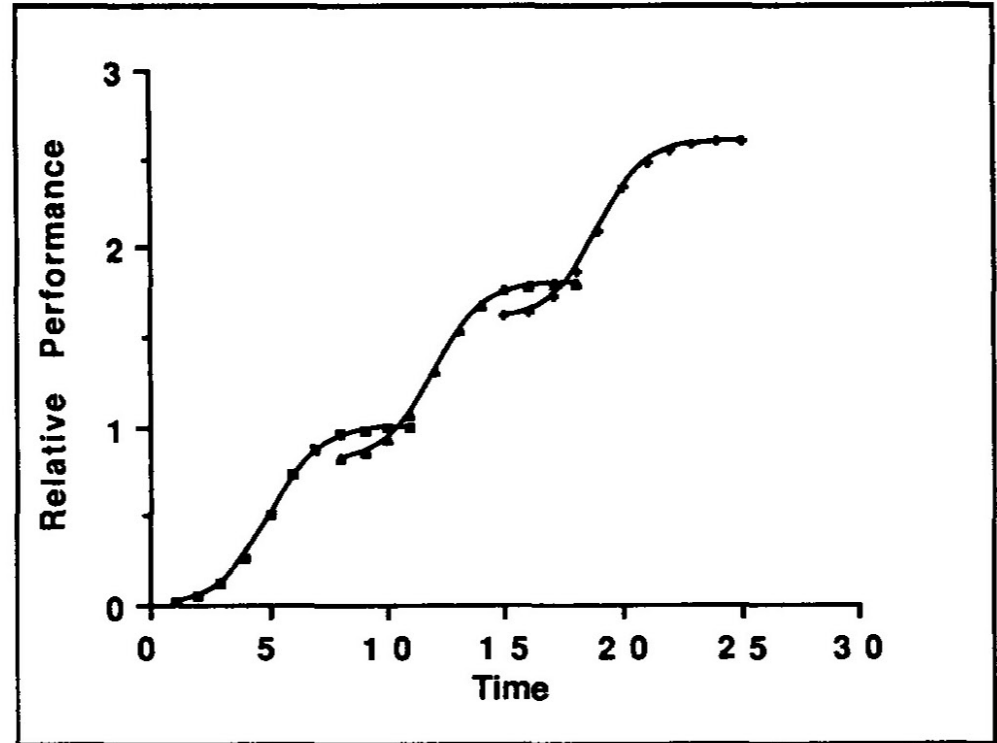
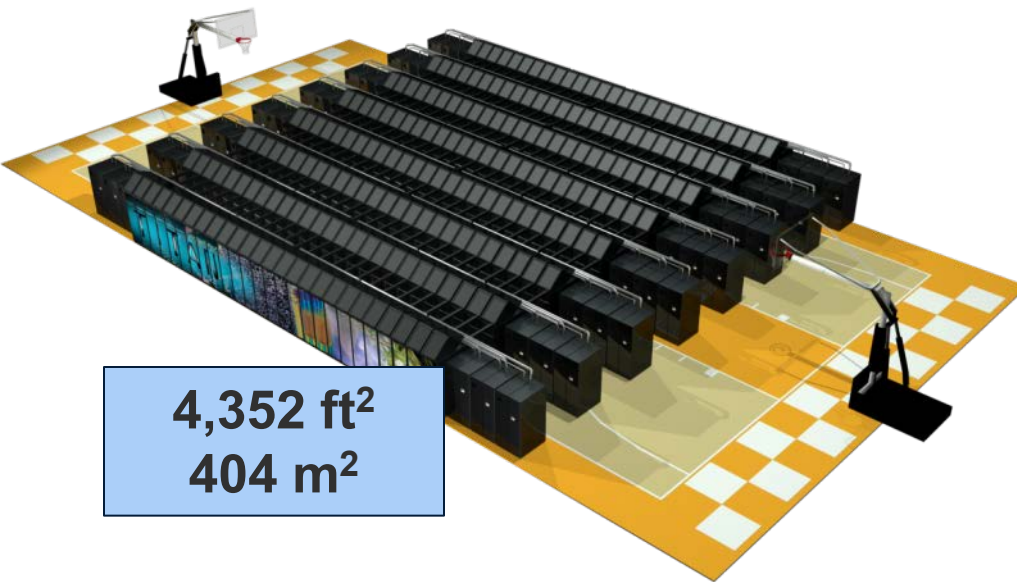
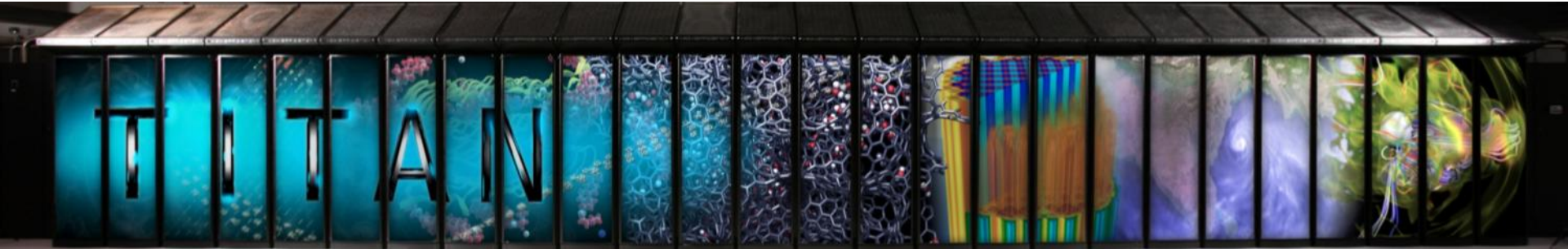


Figure 10. Piecewise-logistic patterns of change.

# The Next Order of Magnitude in Performance

## ORNL's "Titan" Hybrid System: Cray XK7 (AMD x86 Opteron & NVIDIA GPUs)



### SYSTEM SPECIFICATIONS:

- Peak performance of 27.1 PF
  - 24.5 GPU + 2.6 CPU
- 18,688 Compute Nodes each with:
  - 16-Core AMD Opteron CPU
  - NVIDIA Tesla "K20x" GPU
  - 32 + 6 GB memory
- 512 Service and I/O nodes
- 200 Cabinets
- 710 TB total system memory
- Cray Gemini 3D Torus Interconnect
- 8.9 MW peak power



# Cray XK7 Compute Node

## XK7 Compute Node Characteristics

AMD Opteron 6274  
16 core processor @ 141 GF

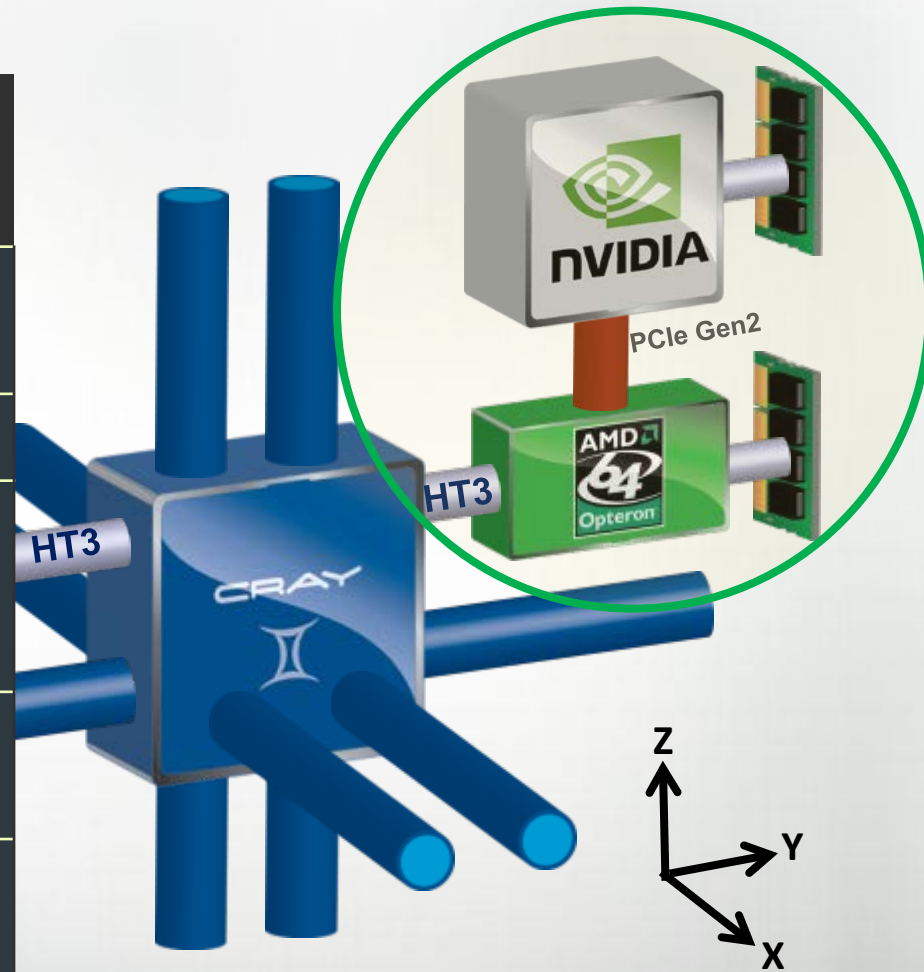
Tesla K20x @ 1311 GF

Host Memory  
32GB

1600 MHz DDR3

Tesla K20x Memory  
6GB GDDR5

Gemini High Speed  
Interconnect

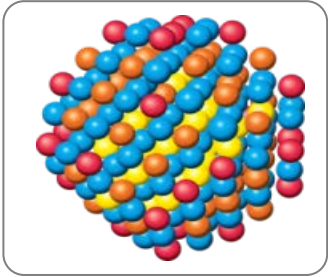


# Center for Accelerated Application Readiness (CAAR)

- **CAAR was created as part of the Titan project to help prepare applications for accelerated architectures**
- **Goals:**
  - **Work with code teams to develop and implement strategies for exposing hierarchical parallelism for our users applications**
  - **Maintain code portability across modern architectures**
  - **Learn from and share our results**
- **Six applications from across different science domains and algorithmic motifs**

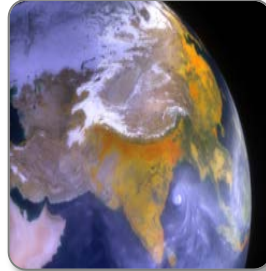


# CAAR Applications on Titan



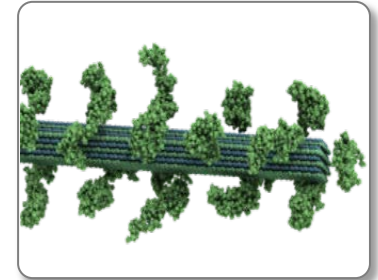
## Material Science (WL-LSMS)

Role of material disorder, statistics, and fluctuations in nanoscale materials and systems.



## Climate Change (CAM-SE)

Answer questions about specific climate change adaptation and mitigation scenarios; realistically represent features like precipitation patterns/statistics and tropical storms.



## Molecular Dynamics (LAMMPS)

A multiple capability molecular dynamics code.

## Astrophysics (NRDF)

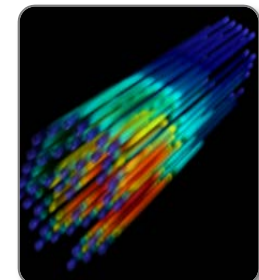
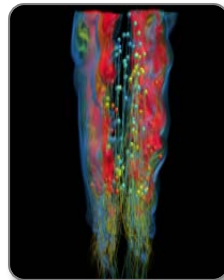
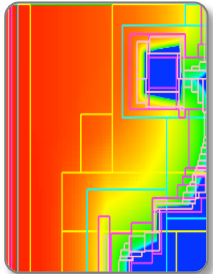
AMR Radiation transport – critical to astrophysics, laser fusion, combustion, atmospheric dynamics, and medical imaging.

## Combustion (S3D)

Combustion simulations to enable the next generation of diesel/bio- fuels to burn more efficiently.

## Nuclear Energy (Denovo)

Unprecedented high-fidelity radiation transport calculations that can be used in a variety of nuclear energy and technology applications.



# How Effective are GPUs on Scalable Applications?

## *OLCF-3 Early Science Codes -- Performance on Titan XK7*

Application	Cray XK7 vs. Cray XE6 Performance Ratio*
<b>LAMMPS*</b> Molecular dynamics	<b>7.4</b>
<b>S3D</b> Turbulent combustion	<b>2.2</b>
<b>Denovo</b> 3D neutron transport for nuclear reactors	<b>3.8</b>
<b>WL-LSMS</b> Statistical mechanics of magnetic materials	<b>3.8</b>
<b>CAM-SE</b> Global Atmospheric Simulation	<b>1.8*</b>

Titan: Cray XK7 (Kepler GPU plus AMD 16-core Opteron CPU)  
Cray XE6: (2x AMD 16-core Opteron CPUs)

\*Performance depends strongly on specific problem size chosen

# Additional Applications from Community Efforts

## *Current Performance Measurements on Titan*

Application	Cray XK7 vs. Cray XE6 Performance Ratio <sup>*</sup>
<b>AWP-ODC</b> Seismology	<b>2.1</b>
<b>DCA++</b> Condensed Matter Physics	<b>4.4</b>
<b>QMCPACK</b> Electronic structure	<b>2.0</b>
<b>RMG (DFT – real-space, multigrid)</b> Electronic Structure	<b>2.0</b>
<b>XGC1</b> Plasma Physics for Fusion Energy R&D	<b>1.8</b>

Titan: Cray XK7 (Kepler GPU plus AMD 16-core Opteron CPU)

Cray XE6: (2x AMD 16-core Opteron CPUs)

<sup>\*</sup>Performance depends strongly on specific problem size chosen



# Some Lessons Learned

- **Exposure of unrealized parallelism**
  - Identifying the opportunities is often straightforward
  - Making changes to exploit it is hard work (made easier by better tools)
  - Developers can quickly learn, e.g., CUDA and put it to effective use
  - A directives-based approach offers a straightforward path to portable performance
- **For those codes that already make effective use of scientific libraries, the possibility of continued use is important.**
  - HW-aware choices
  - Help (or, at least, no hindrance) to overlapping computation with device communication
- **Ensuring that changes are communicated back and remain in the production “trunk” is every bit as important as initially thought**
  - Other development work taking place on all CAAR codes could quickly make acceleration changes obsolete/broken otherwise

# All Codes Will Need Rework To Scale!

- **Up to 2-4 person-years required to port each code from Jaguar to Titan**
  - Refactoring effort will be required for exascale performance regardless of the specific architecture
  - Also pays off for other systems—the ported codes often run significantly faster CPU-only (Denovo 2X, CAM-SE >1.7X)
- **Experience demonstrates 70-80% of developer time is spent in code restructuring, regardless of whether using OpenMP / CUDA / OpenCL / OpenACC / ...**
- **Each code team must make its own choice of using OpenMP vs. CUDA vs. OpenCL vs. OpenACC, based on the specific case—may be different conclusion for each code**
- **The user community and sponsors must plan for this expense**

# All Codes Need Error Recovery at Scale

- **Simple checkpoint / restart is a minimum**
  - At the scale of Titan, we are seeing several nodes fail per day
  - Jobs running on the full system for several hours should expect to have a node fail during job execution and be prepared to recover
- **More advanced error detection and recovery techniques will be required as parallelism increases**
  - FT-MPI, algorithms that can ignore faults, and other research techniques for error containment and recovery mandatory for larger systems

## Need for a richer programming environment

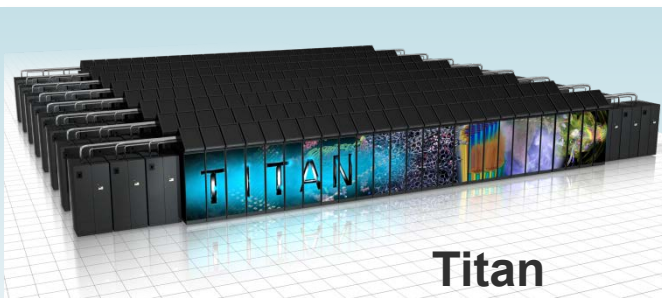
- **Tools are critical to success**
  - complex hierarchical parallelism and heterogeneous processors  $\neq$  the days of debugging with print statements
- **Ongoing investments in good tools essential**
  - debuggers, performance analysis, memory analysis, and the training



# Rethink fundamental algorithmic approach

- **Heterogeneous architectures can make previously intractable or inefficient models and implementations viable**
  - **Alternative methods for electrostatics that perform slower on traditional x86 can be significantly faster on GPUs** (*Nguyen, et al. J. Chem. Theor. Comput. 2013. 73-83*)
  - **3-body coarse-grain simulations of water with greater concurrency can allow > 100X simulation rates when compared to fastest atomistic models even though both are run on the GPUs** (*Brown, et al. Submitted*)

# The Oak Ridge Leadership Computing Facility provides a unique computational user facility for our user community



**Titan**  
**Cray XK7**

Peak performance	27 PF/s
Memory	710 TB
Disk bandwidth	1 TB/s
Square feet	5,000
Power	8.8 MW



**Eos**  
**Cray XC30**

Peak performance	248 TF/s
Memory	48 TB
Disk bandwidth	20 GB/s
Square feet	108

## Data Storage

- Spider File System
  - 40 PB capacity
  - 1+ TB/s bandwidth
- HPSS Archive
  - 240 PB capacity
  - 6 Tape libraries



## Data Analytics & Visualization

- LENS cluster
- Ewok cluster
- Rhea cluster
- EVEREST visualization facility
- uRiKA data appliance



## Networks

- ESnet – 100 Gbps
- Internet2 – 10 Gbps
- XSEDEnet – 10 Gbps
- Private dark fibre

# Summary

- **Partnering has demonstrated value in navigating architectural transition**
  - highly integrated engagement with user community has led to early success
  - CAAR application effort already demonstrating advantages of hybrid architectures
  - user assistance and outreach will help codify best practices and inform the broader community via education and training opportunities
- **Important investments in and collaborations with technology providers**
  - **Scalable Debugging for Hybrid Systems**
    - collaboration with Allinea to develop a scalable hybrid aware debugger based on DDT
  - **High-productivity Hybrid-programming Through Compiler Innovation**
    - collaboration with HMPP to develop directive based compiler technology in CAPS compiler
      - CAPS support for OpenACC set of directives; support for all common languages used at the OLCF
  - **Scalable Performance Analysis for Hybrid Systems**
    - collaboration with Technische Universitat Dresden to add support for Hybrid (CPU/GPU) performance analysis in Vampir



# Questions?



<http://www.nccs.gov>

<http://www.nics.tennessee.edu/>

