#### Exceptional service in the national interest





### Oh, \$#\*@! Exascale!

#### The effect of emerging architectures on scientific discovery Ultrascale Visualization Workshop, November 12, 2012 Kenneth Moreland, Sandia National Laboratories



# Scientific Discovery at the Exascale:

Report from the DOE ASCR 2011 Workshop on Exascale Data Management, Analysis, and Visualization

System Parameter	2011 "2018"		Factor Change	
System Peak	2 PetaFLOPS	1 Exa	FLOP	500
Power	6 MW	≤ 20	MW	3
System Memory	0.3 PB	32 – 64 PB		100 – 200
Total Concurrency	225K	1B × 10	1B × 100	40,000 - 400,000
Node Performance	125 GF	1 TF	10 TF	8 – 80
Node Concurrency	12	1,000	10,000	83 – 830
Network BW	1.5 KB/s	100 GB/s	1000 GB/s	66 – 660
System Size (nodes)	18,700	1,000,000	100,000	50 – 500
I/O Capacity	15 PB	300 – 1000 PB		20 – 67
I/O BW	0.2 TB/s	20 – 60 TB/s		10 - 30

# Scientific Discovery at the Exascale:

Report from the DOE ASCR 2011 Workshop on Exascale Data Management, Analysis, and Visualization

System Parameter	2011	2011 "2018"		Factor Change
System Peak	2 PetaFLOPS	1 ExaFLOP		500
Power	6 MW	≤ 20	MW	3
System Memory	0.3 PB	32 – 64 PB		100 – 200
Total Concurrency	225K	1B × 10	1B × 100	40,000 - 400,000
Node Performance	125 GF	1 TF	10 TF	8 – 80
Node Concurrency	12	1,000	10,000	83 – 830
Network BW	1.5 KB/s	100 GB/s	1000 GB/s	66 – 660
System Size (nodes)	18,700	1,000,000	100,000	50 – 500
I/O Capacity	15 PB	300 – 1000 PB		20 – 67
I/O BW	0.2 TB/s	20 – 60 TB/s		10 - 30

# Scientific Discovery at the Exascale:

Report from the DOE ASCR 2011 Workshop on Exascale Data Management, Analysis, and Visualization

System Parameter	2011 "2018"		18″	Factor Change	
System Peak	2 PetaFLOPS	1 ExaFLOP		500	
Power	6 MW	≤ 20 MW		3	
System Memory	0.3 PB	32 – 64 PB		100 - 200	
Total Concurrency	225K	1B × 10	1B × 100	40,000 – 400,000	
Node Performance	125 GF	1 TF	10 TF	8 – 80	
Node Concurrency	12	1,000	10,000	83 – 830	
Network BW	1.5 KB/s	100 GB/s	1000 GB/s	66 – 660	
System Size (nodes)	18,700	1,000,000	100,000	50 – 500	
I/O Capacity	15 PB	300 – 1	000 PB	20 – 67	
I/O BW	0.2 TB/s	20 – 60	O TB/s	10 - 30	

### **Exascale Projection**



	Jaguar – XT5	Exascale*	Increase
Memory	300 Terabytes	32 – 64 Petabytes	100 – 200×
Concurrency	224,256 way	10 – 100 billion way	Up to 400,000×

MPI Only? Vis object code + state: 20MB On Jaguar: 20MB × 200,000 processes = 4TB On Exascale: 20MB × 100 billion processes = 2EB !

\*Source: Scientific Discovery at the Exascale, Ahern, Shoshani, Ma, et al.

### **Exascale Projection**



	Jaguar – XT5	Exascale*	Increase
Memory	300 Terabytes	32 – 64 Petabytes	100 – 200×
Concurrency	224,256 way	10 – 100 billion way	Up to 400,000×

Visualization pipeline too heavyweight? On Jaguar: 1 trillion cells  $\rightarrow$  5 million cells/thread On Exascale: 100 trillion cells  $\rightarrow$  1000 cells/thread

\*Source: Scientific Discovery at the Exascale, Ahern, Shoshani, Ma, et al.

### **Exascale Projection**



	Jaguar – XT5	Exascale*	Increase
Memory	300 Terabytes	32 – 64 Petabytes	100 – 200×
Concurrency	224,256 way	10 – 100 billion way	Up to 400,000×

### Overhead of ghost/halo cells?

On Jaguar: 1 trillion cells → 5 million cells/thread
Partition into ~171<sup>3</sup> blocks
6 × 171<sup>2</sup> ≈ 175K ghost/block → 35 billion ghost total
Ghost cells ~3.5% size of original data
On Exascale: 100 trillion cells → 1000 cells/thread
Partition into 10<sup>3</sup> blocks
6 × 10<sup>2</sup> ≈ 600 ghost/block → 60 trillion ghost total
Ghost cells 60% size of original data

\*Source: Scientific Discovery at the Exascale, Ahern, Shoshani, Ma, et al.

#### SDAV PA

## Exascale Programming Challenges

Sandia National Laboratories

- At some point, domain decomposition fails
  - Too many halo cells, too much communication
- Possible new architectures and programming models
  - GPU accelerators hate decomposition
- Threaded (OpenMP) programming is easier than distributed (MPI) programming.
  - Threading needs careful planning for memory affinity (inherent in distributed)
  - Sharing memory locations invites read/write collisions (explicit in distributed)
  - PGAS will save us? I'm skeptical.
- Best practice approach: Encapsulated Multithreaded Operations
  - Multiple DOE projects underway: Dax (ASCR), PISTON (ASC), EAVL (LDRD)
  - If successful, minimal impact on applications
  - Might be some changes in scope of what can be done

# Scientific Discovery at the Exascale:

Report from the DOE ASCR 2011 Workshop on Exascale Data Management, Analysis, and Visualization

System Parameter	2011 "2018"		Factor Change	
System Peak	2 PetaFLOPS	1 ExaFLOP		500
Power	6 MW	≤ 20	MW	3
System Memory	0.3 PB	32 – 64 PB		100 – 200
Total Concurrency	225K	1B × 10	1B × 100	40,000 - 400,000
Node Performance	125 GF	1 TF	10 TF	8 – 80
Node Concurrency	12	1,000	10,000	83 – 830
Network BW	1.5 KB/s	100 GB/s	1000 GB/s	66 – 660
System Size (nodes)	18,700	1,000,000	100,000	50 – 500
I/O Capacity	15 PB	300 – 1000 PB		20 – 67
I/O BW	0.2 TB/s	20 – 60 TB/s		10 - 30

## Extreme scale computing

SDAV DAV Sandia National

- Trends
  - More FLOPS
  - More concurrency
  - Comparatively less storage,
     I/O bandwidth
- ASCI purple (49 TB/140 GB/s) – JaguarPF (300 TB/ 200 GB/s)
  - Most people get < 5 GB/sec at scale





Computation 1 EB/s

Node Memory 400 PB/s

Interconnect (10% Staging Nodes) 10 PB/s

Storage 60 TB/s

Computation 1 EB/s

Node Memory 400 PB/s

Interconnect (10% Staging Nodes) 10 PB/s

Off-Line Visualization

Storage 60 TB/s Embedded Visualization Computation 1 EB/s

Node Memory 400 PB/s

Co-Scheduled Visualization

Interconnect (10% Staging Nodes) 10 PB/s

Off-Line Visualization

Storage 60 TB/s

### Space of Solutions



	Capability	Coupling	Footprint	Transfer	Interactive
Tightly Integrated	Low	Tight	Low	None	No
Embedded	High	Tight	High	Possible memcpy	No
Hybrid	High	Tight	Medium	Subset Hi Speed Transfer	Yes
Co-Scheduled	High	Loose	~5% Extra Nodes	Hi Speed Transfer	Yes
Off-Line	High	Loose	None	Slow Persistent Storage Cost	Yes

### **Other Exascale Challenges**



#### Resilience

- Mean time to failure
- Robustness for in situ
- Soft errors
- Uncertainty
- Compression/extraction
  - Feature characterization
  - Lossy/lossless compression
- Provenance
  - Capture functions, algorithms, parameters
- Uncertainty Quantification

## Acknowledgements



- Funding support
  - The DOE Office of Science, Advanced Scientific Computing Research, under award number 10-014707, program manager Lucy Nowell
  - The Director, Office of Advanced Scientific Computing Research, Office of Science, of the U.S. Department of Energy under Contract No. 12-015215, through the Scientific Discovery through Advanced Computing (SciDAC) Institute of Scalable Data Management, Analysis and Visualization
- Collaborators
  - Sandia National Laboratories: Kenneth Moreland, Nathan Fabian, Ron Oldfield
  - Los Alamos National Laboratory: James Ahrens, Jonathan Woodring
  - Oak Ridge National Laboratory: Scott Klasky, Norbert Podhorszki
  - Argonne National Laboratory: Venkatram Vishwanath, Mark Hereld, Michael E.
     Papka
  - Kitware, Inc.: Berk Geveci, Utkarsh Ayachit, Andrew C. Bauer, Pat Marion, Sebastien Jourdain, David DeMarle, David Thompson
  - University of Colorado at Boulder: Michel Rasquin, Kenneth E. Jansen
  - Rutgers University: Ciprian Docan, Manish Parashar