

NCAR's Data-Centric Supercomputing Environment Yellowstone

November 29, 2011
David L. Hart, CISL
dhart@ucar.edu



Welcome to the Petascale

- **Yellowstone hardware and software**
- **Deployment schedule**
- **Allocations opportunities at NWSC**
 - ASD, University, CSL, NCAR, and Wyoming-NCAR alliance

Construction complete!



NCAR Resources

at the NCAR-Wyoming Supercomputing Center (NWSC)

- **Centralized Filesystems and Data Storage**

- 10.9 PB initially → 16.4 PB in 1Q2014
- 12x usable capacity of current GLADE

- **High-Performance Computing**

- IBM iDataPlex cluster
- 1.55 PFLOPs
- 30x Bluefire computing performance

- **Data Analysis and Visualization**

- Large-memory system
- GPU computation and visualization system
- Knights Corner system

NWSC-1
Procurement

- **NCAR HPSS Data Archive**

- 2 StorageTek SL8500 tape libraries – 20k cartridge slots
- >100 PB capacity with 5 TB cartridges (uncompressed)

AMSTAR
Procurement

Arrived
Nov. 4

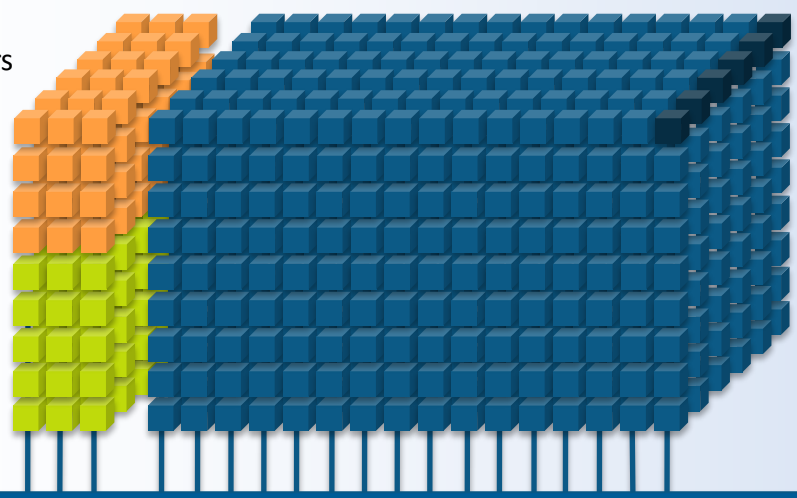
Yellowstone Environment

Computational & Information Systems Laboratory

Geyser & Caldera
DAV clusters

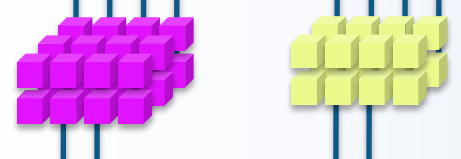
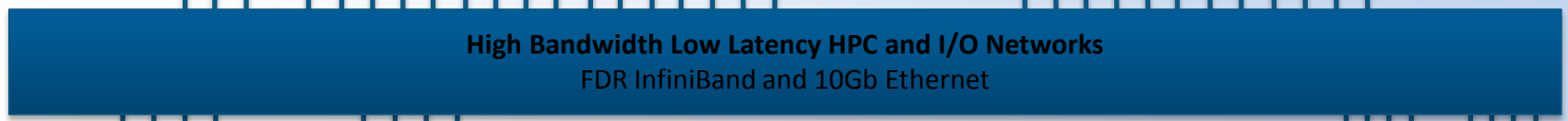
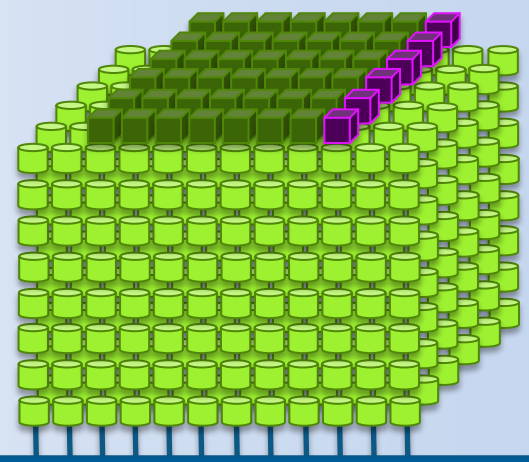
Yellowstone

HPC resource, 1.55 PFLOPS peak



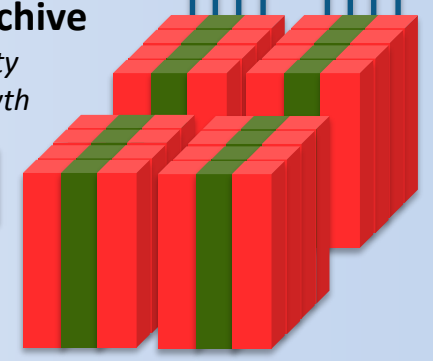
GLADE

Central disk resource
11 PB (2012), 16.4 PB (2014)



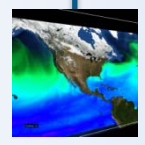
NCAR HPSS Archive

100 PB capacity
~15 PB/yr growth



Science Gateways
RDA, ESG

Data Transfer
Services



Remote Vis

Partner Sites

XSEDE Sites



Yellowstone

NWSC High-Performance Computing Resource

- **Batch Computation**

- 74,592 cores total – 1.552 PFLOPs peak
- 4,662 IBM dx360 M4 nodes – 16 cores, 32 GB memory per node
- Intel Sandy Bridge EP processors with AVX – 2.6 GHz clock
- 149.2 TB total DDR3-1600 memory
- 29.8 Bluefire equivalents

- **High-Performance Interconnect**

- Mellanox FDR InfiniBand full fat-tree
- 13.6 GB/s bidirectional bw/node
- <2.5 μ s latency (worst case)
- 31.7 TB/s bisection bandwidth

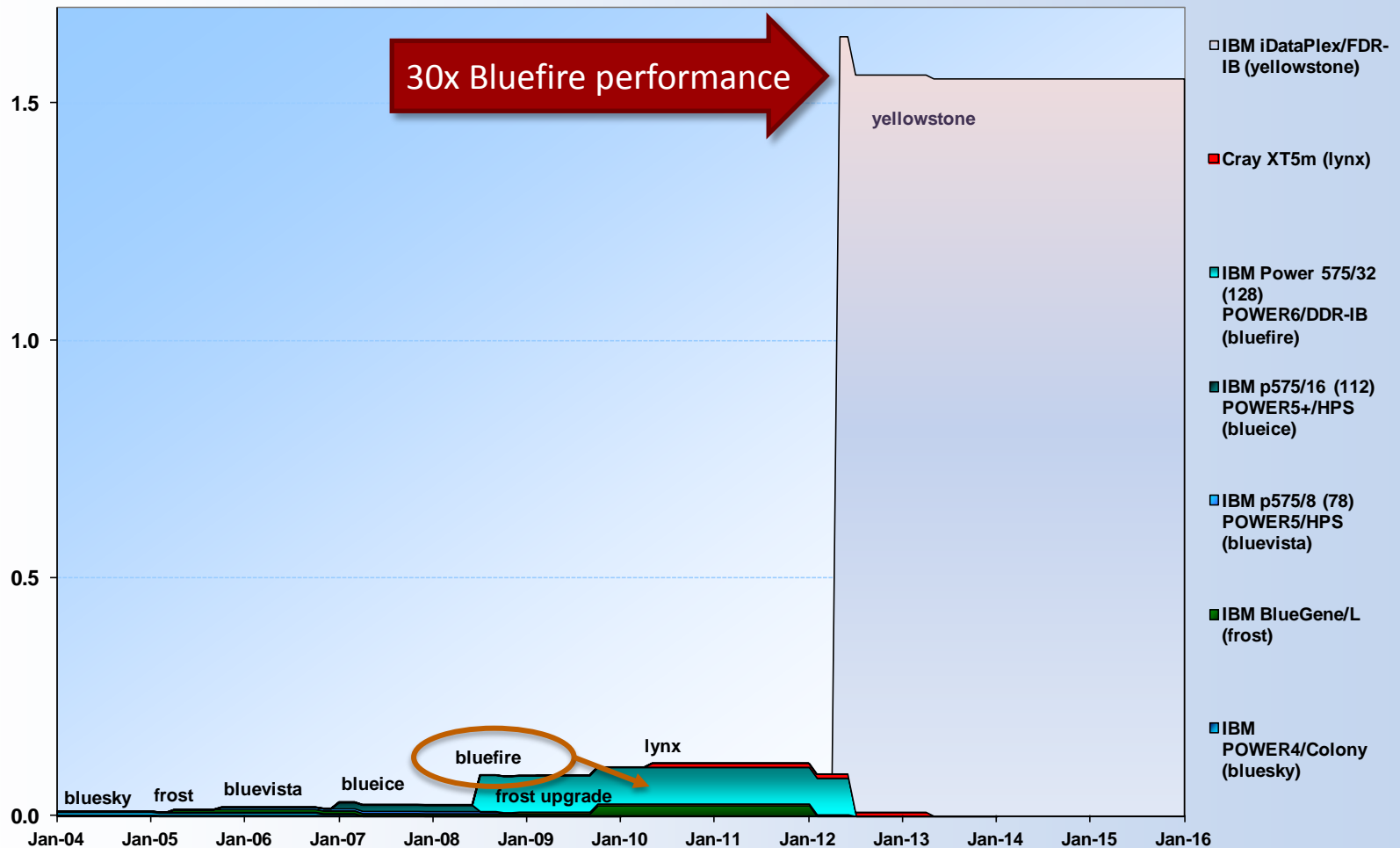
- **Login/Interactive**

- 6 IBM x3650 M4 Nodes; Intel Sandy Bridge EP processors with AVX
- 16 cores & 128 GB memory per node



NCAR HPC Profile

Peak PFLOPs at NCAR



GLADE

- **10.94 PB usable capacity → 16.42 PB usable (1Q2014)**

Estimated initial file system sizes

- **collections** ≈ 2 PB RDA, CMIP5 data
- **scratch** ≈ 5 PB shared, temporary space
- **projects** ≈ 3 PB long-term, allocated space
- **users** ≈ 1 PB medium-term work space

- **Disk Storage Subsystem**

- 76 IBM DCS3700 controllers & expansion drawers
 - 90 2-TB NL-SAS drives/controller
 - add 30 3-TB NL-SAS drives/controller (1Q2014)

- **GPFS NSD Servers**

- **91.8 GB/s** aggregate I/O bandwidth; 19 IBM x3650 M4 nodes

- **I/O Aggregator Servers (GPFS, GLADE-HPSS connectivity)**

- 10-GbE & FDR interfaces; 4 IBM x3650 M4 nodes

- **High-performance I/O interconnect to HPC & DAV**

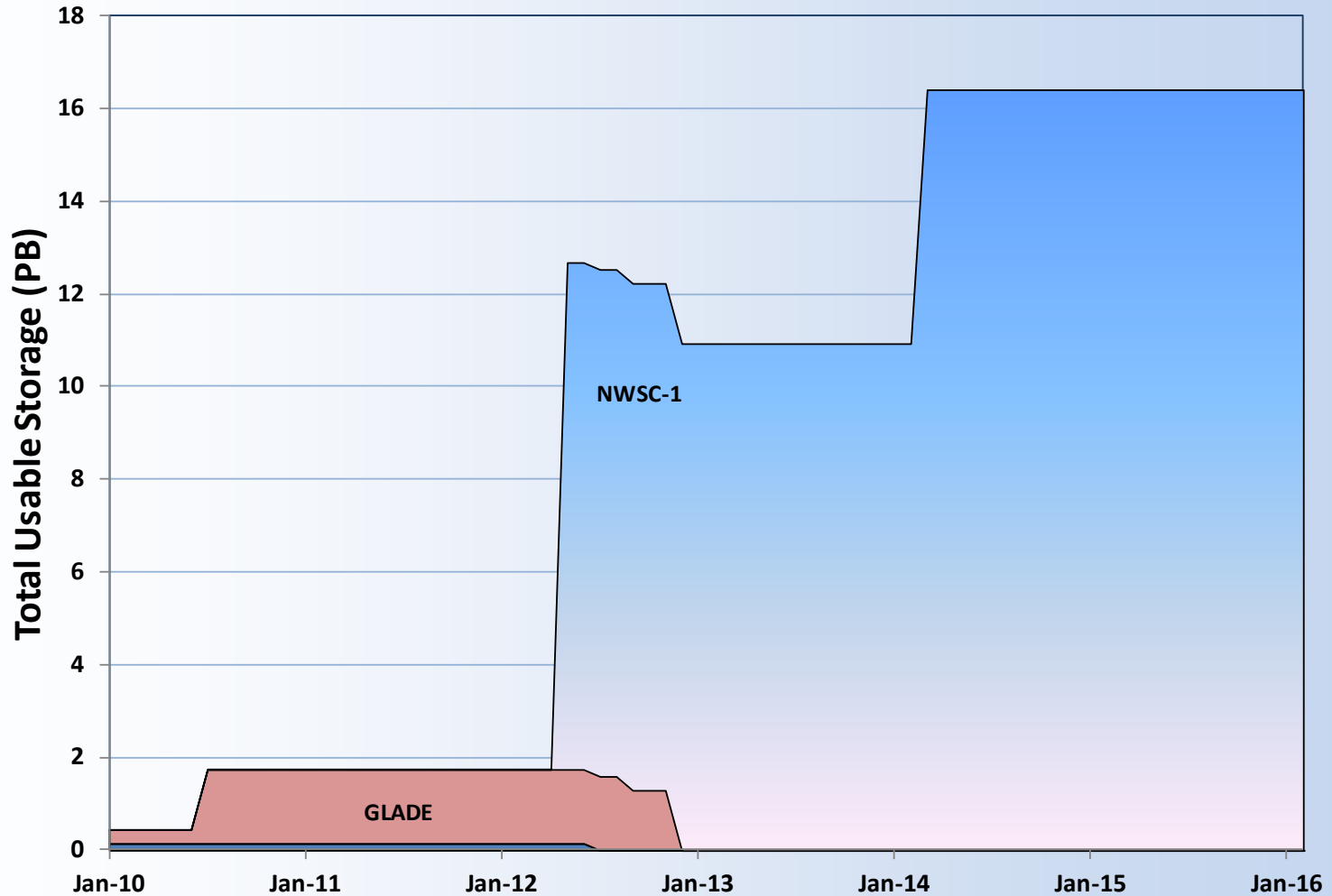
- Mellanox FDR InfiniBand full fat-tree
- 13.6 GB/s bidirectional bandwidth/node



NCAR Disk Capacity Profile

Total Centralized Filesystem Storage (PB)

NWSC-1 GLADE bluefire



Geyser and Caldera

NWSC Data Analysis & Visualization Resource

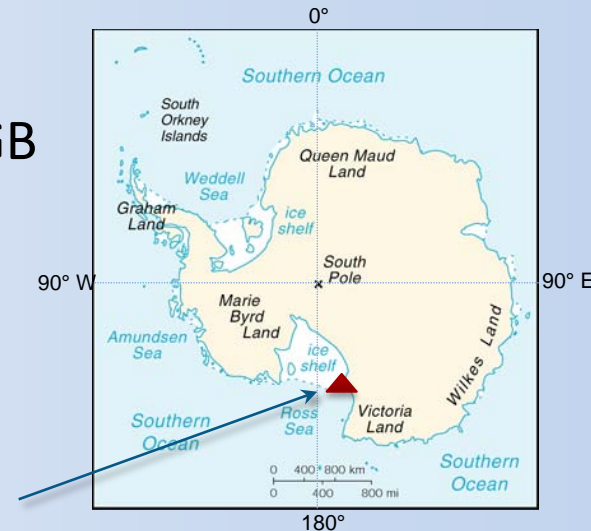
- **Geyser: Large-memory system**
 - 16 IBM x3850 nodes – Intel Westmere-EX processors
 - 40 cores, **1 TB memory**, 1 NVIDIA Kepler Q13H-3 GPU *per node*
 - Mellanox FDR full fat-tree interconnect
- **Caldera: GPU computation/visualization system**
 - 16 IBM x360 M4 nodes – Intel Sandy Bridge EP/AVX
 - 16 cores, 64 GB memory per node
 - 2 NVIDIA Kepler Q13H-3 GPUs per node
 - Mellanox FDR full fat-tree interconnect
- **Knights Corner system (November 2012 delivery)**
 - Intel Many Integrated Core (MIC) architecture
 - 16 IBM Knights Corner nodes
 - 16 Sandy Bridge EP/AVX cores, 64 GB memory
 - 1 Knights Corner adapter per node
 - Mellanox FDR full fat-tree interconnect



Erebus

Antarctic Mesoscale Prediction System (AMPS)

- **IBM iDataPlex Compute Cluster**
 - 84 IBM dx360 M4 Nodes; 16 cores, 32 GB
 - Intel Sandy Bridge EP; 2.6 GHz clock
 - 1,344 cores total – 28 TFLOPs peak
 - Mellanox FDR InfiniBand full fat-tree
 - 0.54 Bluefire equivalents
- **Login Nodes**
 - 2 IBM x3650 M4 Nodes
 - 16 cores & 128 GB memory per node
- **Dedicated GPFS filesystem**
 - 57.6 TB usable disk storage
 - 9.6 GB/sec aggregate I/O bandwidth



Erebus, on Ross Island, is Antarctica's most famous volcanic peak and is one of the largest volcanoes in the world – within the top 20 in total size and reaching a height of 12,450 feet.



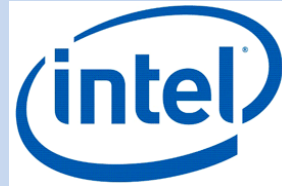
Yellowstone Software

- **Compilers, Libraries, Debugger & Performance Tools**

- **Intel** Cluster Studio (Fortran, C++, performance & MPI libraries, trace collector & analyzer) 50 concurrent users
- **Intel** VTune Amplifier XE performance optimizer 2 concurrent users
- **PGI** CDK (Fortran, C, C++, pgdbg debugger, pgprof) 50 conc. users
- **PGI** CDK GPU Version (Fortran, C, C++, pgdbg debugger, pgprof) for DAV systems only, 2 concurrent users
- **PathScale** EccoPath (Fortran C, C++, PathDB debugger) 20 concurrent users
- Rogue Wave **TotalView** debugger 8,192 floating tokens
- **IBM** Parallel Environment (POE), including IBM HPC Toolkit

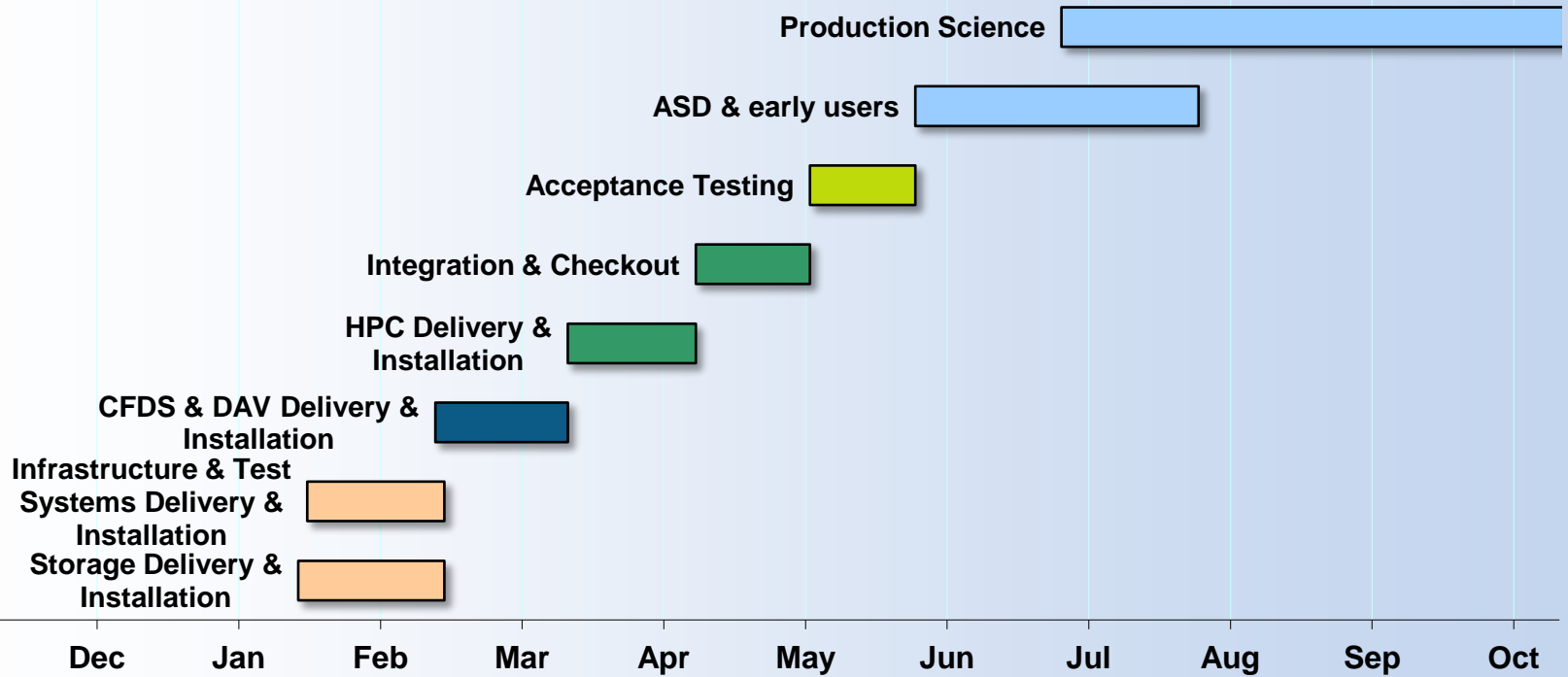
- **System Software**

- **LSF-HPC** Batch Subsystem / Resource Manager
 - IBM has purchased Platform Computing, Inc., developers of LSF-HPC
- Red Hat Enterprise **Linux** (RHEL) Version 6
- IBM General Parallel Filesystem (**GPFS**)
- Mellanox Universal Fabric Manager
- IBM xCAT cluster administration toolkit



Yellowstone Schedule

Delivery, Installation, Acceptance & Production



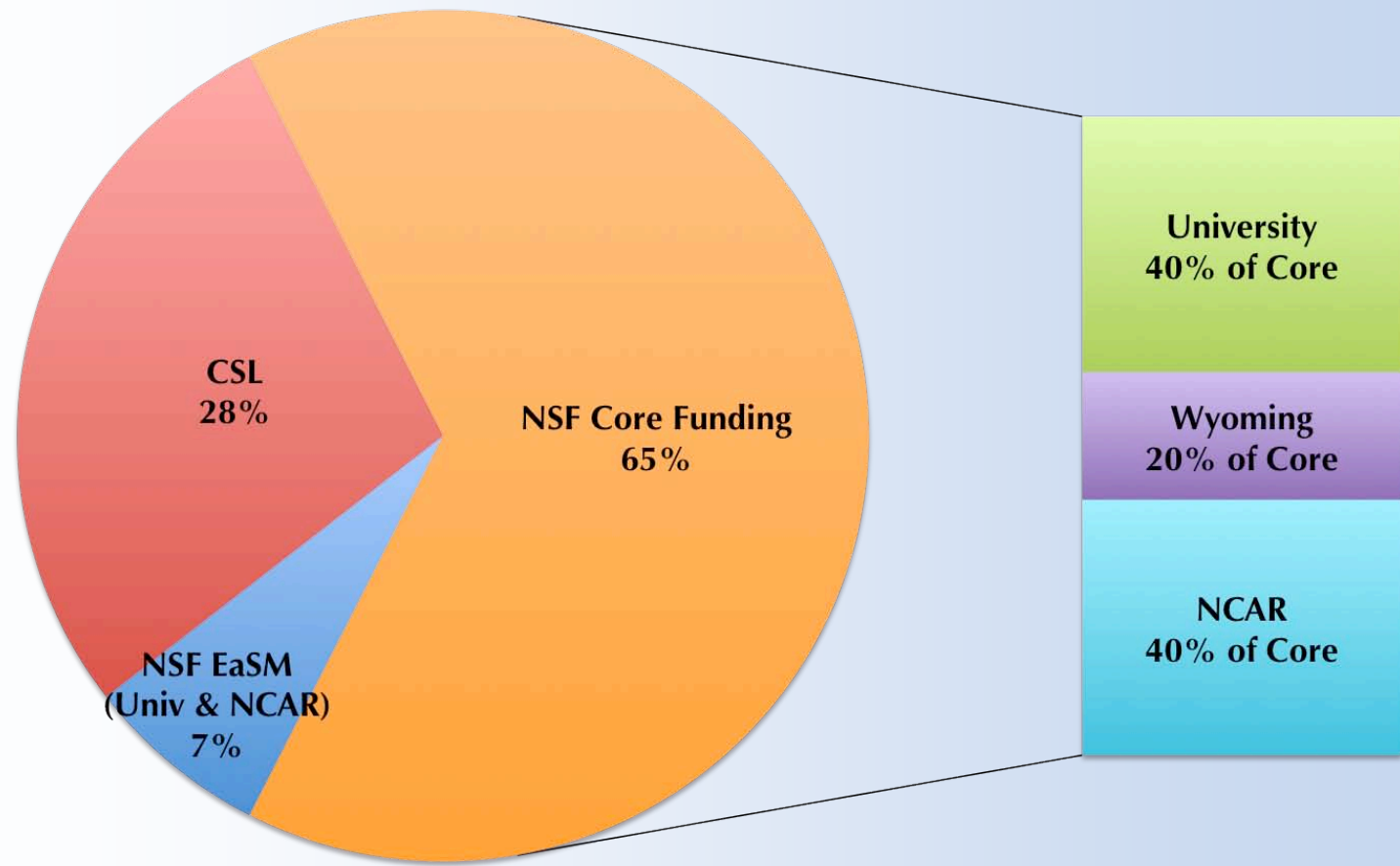
Janus: Available now

- **Janus Dell Linux cluster**
 - 16,416 cores total – 184 TFLOPs peak
 - 1,368 nodes – 12 cores, 24 GB memory per node
 - Intel Westmere processors – 2.8 GHz clock
 - 32.8 TB total memory
 - QDR InfiniBand interconnect
 - Red Hat Linux, Intel compilers (PGI coming)
- **Deployed by CU in collaboration with NCAR**
 - ~10% of the system allocated by NCAR
- ***Available for Small allocations to university, NCAR users***
 - CESM, WRF already ported and running
 - Key elements of NCAR software stack already installed
- **www2.cisl.ucar.edu/docs/janus-cluster**



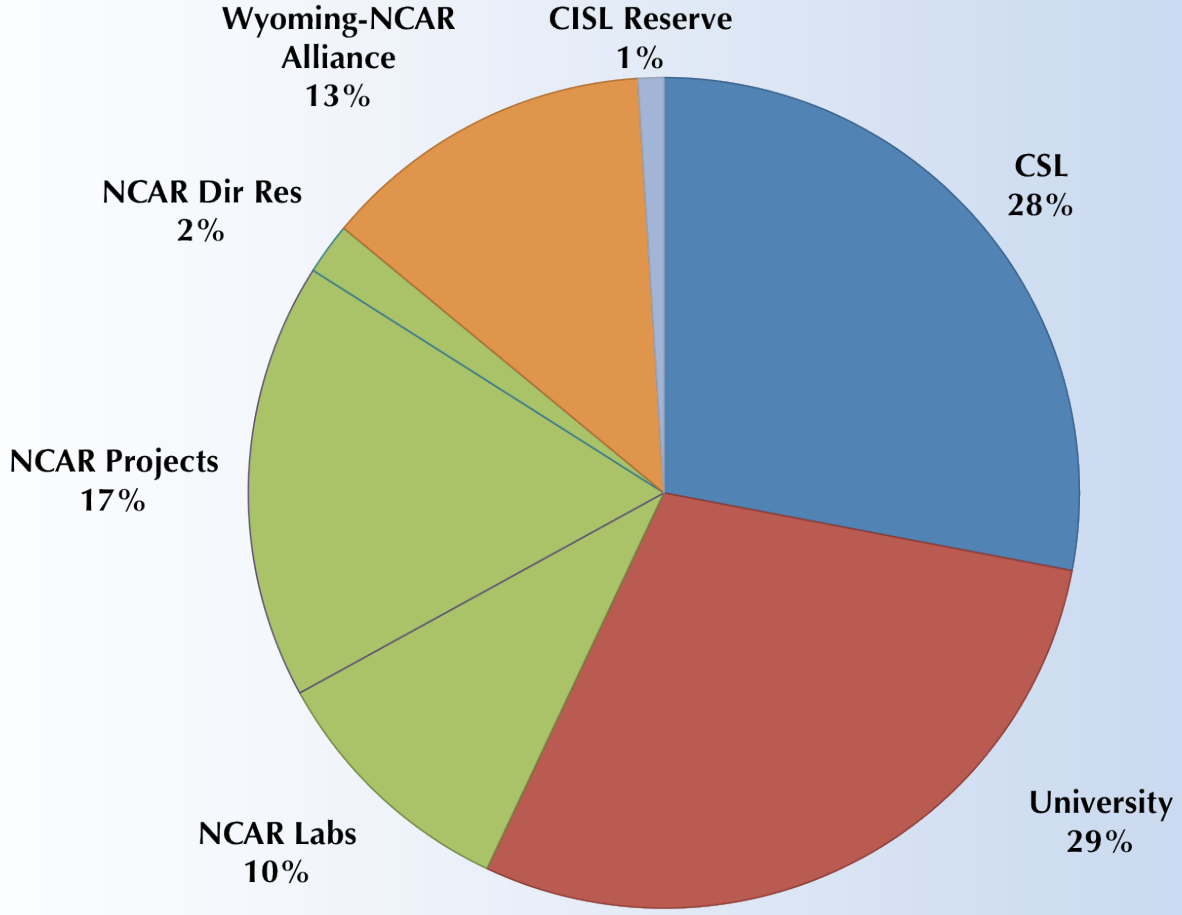


Yellowstone allocation opportunities



Yellowstone funding sources

Yellowstone will be capable of *653 million core-hours per year*, compared to 34 million for Bluefire, and each Yellowstone core-hour is equivalent to 1.53 Bluefire core-hours.



Yellowstone allocations opportunities

The segments for CSL, University and NCAR users each represent about *170 million core-hours per year* on Yellowstone (compared to less than 10 million per year on Bluefire) plus a similar portion of DAV and GLADE resources.

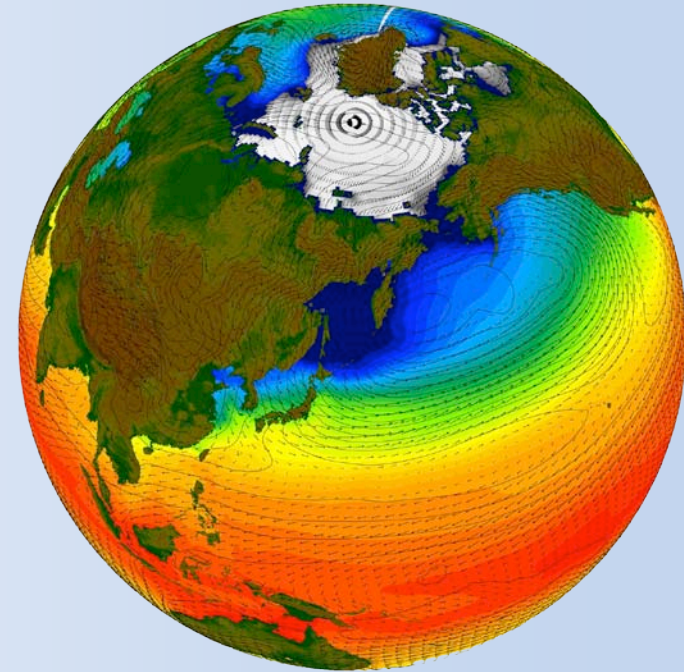
Early-use opportunity:

Accelerated Scientific Discovery

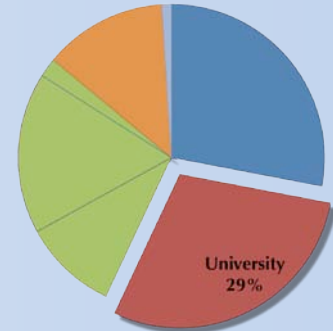
- **Deadline: January 13, 2012**
- **Targeting a small number of rapid-turnaround, large-scale projects**
 - *Minimum* HPC request of 5 million core hours
 - Roughly May-July, with access to DAV systems beyond that point through final report deadline, February 2013
- **Approximately 140 million core-hours, in two parts**
 - University-led projects with NSF awards in the geosciences will be allocated **70 million core-hours**
 - NCAR-led projects will make up the other half, selected from NCAR Strategic Capability requests that designate themselves “ASD-ready.”
- **Particularly looking for projects that contribute to NWSC Community Science Objectives**
 - High bar for production readiness, including availability of staff time
- **www2.cisl.ucar.edu/docs/allocations/asd**

Climate Simulation Laboratory

- ***Deadline: mid-February 2012***
- **Targets large-scale, long-running simulations of the Earth's climate**
 - Dedicated facility supported by the U.S. Global Change Research Program
 - Must be climate-related work, but support may be from any agency
- ***Minimum request and award size***
 - Typically 18-month allocation period
 - Approx. 250 million core-hours to be allocated
 - Estimated minimum request size: ~10 million core-hours
- **Preference given to large, collective group efforts, preferably interdisciplinary teams**



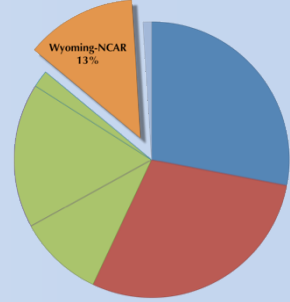
University allocations



- ***Next deadline: March 26, 2012***
- **Large allocations will continue to be reviewed and awarded twice per year**
 - Deadlines in March and September
 - Approx. 85 million core-hours to be allocated at each opportunity
- **Small allocations will also be available once system enters full production**
 - “Small” threshold still to be determined
 - Small allocations for researchers with NSF award—appropriate for benchmarking, preparation for large request
 - Small, one-time allocations for grad students, post-docs, new faculty without NSF award
 - Classroom allocations for instructional use
- ***www2.cisl.ucar.edu/docs/allocations/university***

NEW!

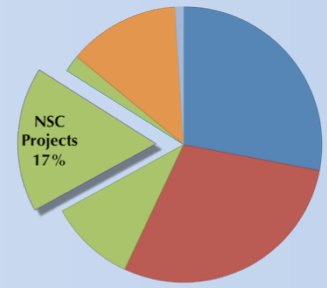
Wyoming-NCAR Alliance



- **Deadline: March 2012 (TBD)**
- **13% of Yellowstone resources**
 - 75 million core-hours per year
 - U Wyoming managed process
- **Activities must have substantial U Wyoming involvement**
 - Allocated projects must have Wyoming lead
 - Extended list of eligible fields of science
 - Eligible funding sources not limited to NSF
- **Actively seeking to increase collaborations with NCAR and with other EPSCoR states.**
- **Otherwise, process modeled on University allocations, with panel review of large requests.**
- ***www.uwyo.edu/nwsc***

NEW!

NCAR Strategic Capability projects



- **First submission deadline: January 13, 2012**
- **17% of Yellowstone resources**
 - 100 million core-hours per year
 - 80x the NCAR Capability Computing annual levels
- **NCAR-led activities linked to NCAR scientific/strategic priorities and/or NWSC science justification**
 - Target large-scale projects within finite time period (one year to a few years)
 - Minimum of 5 million core-hours (exceptions possible)
 - Proposed and reviewed each year
- **Submission format and review criteria similar to that used for large University requests**
- **NCAR ASD projects will be selected from requests here that designate themselves “ASD-ready”**
- ***www2.cisl.ucar.edu/docs/allocations/ncar***

Allocation changes in store

- **Not just HPC, but DAV, HPSS, GLADE allocations**
 - Non-HPC resources \approx 1/3 procurement cost
 - Ensure that use of scarce and costly resources are directed to the most meritorious projects
- **Balance between the time to prepare and review requests and the resources provided**
 - Minimize user hurdles and reviewer burden
 - Build on familiar process for requesting HPC allocations
- **Want to identify projects contributing to the NWSC Community Scientific Objectives**
 - *www2.cisl.ucar.edu/resources/yellowstone/science*
- **All new, redesigned accounting system (SAM)**
 - Separate, easier to understand allocations
 - Switchable 30/90 option, per project, as an operational control
 - (“30/90” familiar to NCAR labs and CSL awardees)

General submission format

- ***Please see specific opportunities for detailed guidelines!***
- **Five-page request**
 - A. Project information (title, lead, etc.)
 - B. Project overview and strategic linkages
 - C. Science objectives
 - D. Computational experiments and resource requirements (HPC, DAV, and storage)
- **Supporting information**
 - E. Multi-year plan (if applicable)
 - F. Data management plan
 - G. Accomplishment report
 - H. References and additional figures



Tips and advice

- **Remember your audience: computational geoscientists from national labs, universities and NCAR**
 - Don't assume they are experts in *your* specialty
- **Be sure to articulate relevance and linkages**
 - Between funding award, computing project, eligibility criteria, and NWSC science objectives (as appropriate)
- ***Don't submit a science proposal***
 - Describe the science in detail sufficient to justify the computational experiments proposed
- **Most of the request should focus on computational experiments and resource needs**
 - Effective methodology
 - Appropriateness of experiments
 - Efficiency of resource use

Justifying resource needs

- **HPC — similar to current practice**
 - Cost of runs necessary to carry out experiment, supported by benchmark runs or published data
- **DAV — will be allocated, similar to HPC practice**
 - A “small” allocation will be granted upon request
 - Allocation review to focus on larger needs associated with batch use
 - Memory and GPU charging to be considered
- **HPSS — focus on storage needs above a threshold**
 - 20-TB default threshold initially
 - Perhaps lower default for “small” allocations”
 - CISL to evaluate threshold regularly to balance requester/reviewer burden with demand on resources
 - Simplified request/charging formula
- **GLADE — project (long-term) spaces will be reviewed and allocated**
 - scratch, user spaces not allocated

GLADE resource requests

- **Only for project space**
 - No need to detail use of scratch, user spaces
- **Describe why project space is essential**
 - That is, why scratch or user space insufficient
 - Show that you are aware of the differences
 - Shared data, frequently used, not available on disk from RDA, ESG (collections space)
- **Relate the storage use to your workflow and computational plan**
 - Projects with data-intensive workflows should show they are using resources efficiently

HPSS resource requests

- **Goal: Demonstrate that HPSS use is efficient and appropriate**
 - Not fire and forget into a “data coffin”
 - Not using as a temporary file system
- **Explain new data to be generated**
 - Relate to computational experiments proposed
 - Describe scientific value/need for data stored
- **Justify existing stored data**
 - Reasons for keeping, timeline for deletion
- **Data management plan: Supplementary information**
 - Additional details on the plans and intents for sharing, managing, analyzing, holding the data

<http://www2.cisl.ucar.edu/resources/yellowstone>
<http://www2.cisl.ucar.edu/docs/allocations>
cislhelp@ucar.edu or dhart@ucar.edu



QUESTIONS?