

# Requirements for Next-Generation OLCF Systems

Lattice QCD Computational Science Workshop  
April 29-30, 2013



Wayne Joubert

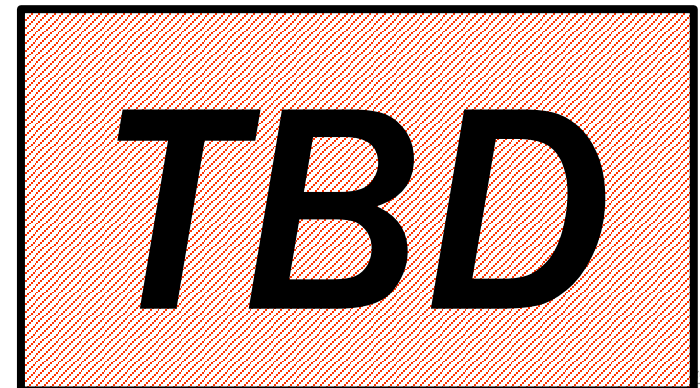
Scientific Computing Group

Oak Ridge Leadership Computing Facility

Oak Ridge National Laboratory

# Context: OLCF-4

- We are now in the planning phase for our next system, a pre-exascale system of 100-200 PF capability which will be deployed in the 2016-2017 time frame (OLCF-4)
- We are currently in a requirements gathering process with users and projects, to better understand what they will require of this system and how they plan to use it



# The Requirements Process

- **Requirement:**
  - “a condition or capability needed by a user to solve a problem or achieve an objective”
- **Requirements for OLCF Leadership Computing:**
  - Items/specifications needed in order for teams to meet science objectives of OLCF projects and the DOE
- **Example areas:**

system hardware requirements	science model requirements
system software requirements	algorithm requirements
compilers, libraries, tools	application codes
data storage, access, analysis	development processes
workflow management	

# Requirements for Leadership Computing Facilities

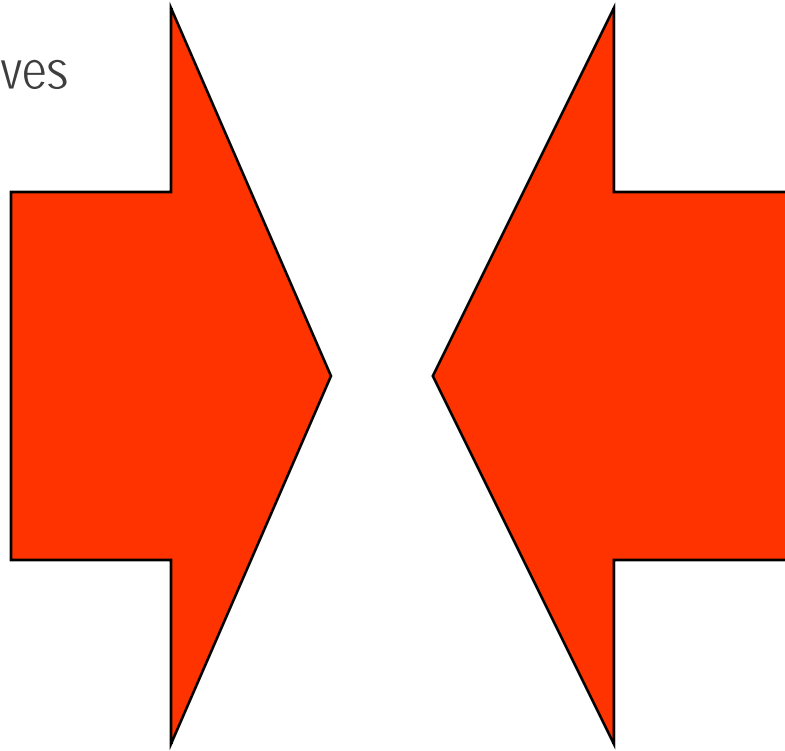
- ORNL leadership computing systems focus on “capability” computing rather than “capacity” computing—to solve problems that can’t be solved elsewhere
- Future requirements estimation in this situation is a challenge—our customers and workloads can vary from year to year based on what projects are awarded time
- We must estimate required future machine characteristics based on current projects, assuming slow change in the mix of projects



# Requirements Analysis: A 2-Way Street

“What kind of science can be done”

- ✓ Mission, project objectives
- ✓ New problem regimes
- ✓ New models
- ✓ New algorithms
- ✓ Higher accuracy
- ✓ More DOF
- ✓ More timesteps
- ✓ Faster time to solution
- ✓ In-depth analytics
- ✓ UQ



“What kind of system can be built”

- ✓ Moore’s law
- ✓ Power constraints
- ✓ Budget constraints
- ✓ Mass-market hardware trends
- ✓ Processor trends – e.g., GPUs
- ✓ New interconnects
- ✓ Storage

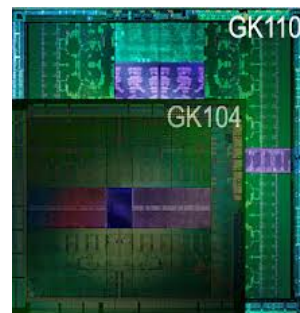
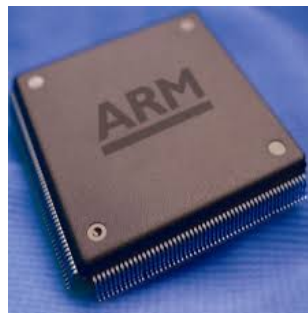
**Code developers must understand system capabilities and limitations, and HPC architects must address the needs of a wide range of science domain codes.**

# Meeting Mission Requirements in a Changing HPC Environment

To meet science goals, we are being required to embrace increasingly disruptive changes in HPC hardware and software—with more to come

## ➤ The hardware complexity challenge

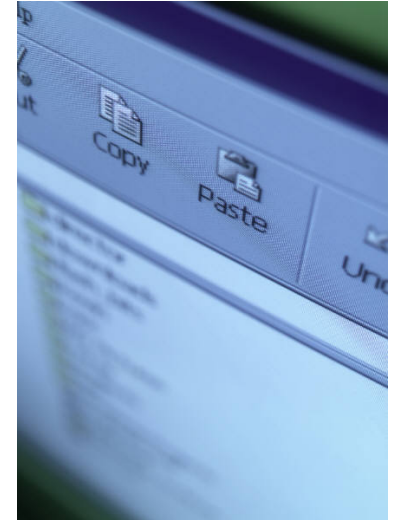
- Higher chip densities, more complex nodes, heterogeneous processors, processor diversity – “compute jungle”
- More difficult to program, potentially harder to diagnose faults and performance issues
- Increasing resilience concern as failure rates trend upward



# The Changing HPC Environment

## ➤ The software complexity challenge

- Software stack is growing in complexity to support more complex hardware
- Tendency of application codes, models, algorithms to grow in complexity



# The Changing HPC Environment

## ➤ The programming model challenge

- Presently to get good performance an application may need to support MPI, OpenMP, CUDA or vectorization directives or function calls, processor/memory heterogeneity, checkpoint restart.
- In the near future, we may need NVRAM out-of-core programming, user-managed fault tolerance, optimizing code for interconnect topologies, energy-aware programming and possibly programming to custom hardware.
- The user-facing programming environment must be improved.

## ➤ The algorithm challenge

- Hardware changes may require new kinds of algorithms, e.g., out-of-core, communication-avoidant algorithms, cache-aware/cache-oblivious techniques, parallelism in time, other potentially revolutionary algorithm changes



# The OLCF Requirements Process

We collect information from our users regarding concerns and needs such as these

We collect information from vendors about what is possible for next-generation systems

We feed this information into our planning process

# Example: The 2009 OLCF Requirements Process

During the 2009 requirements process, we discovered:

- “more flops” were needed broadly across many projects
- Some codes were starting to reach limits to the extent to which they could “scale out” (e.g., strong scaling), now needed to “scale in” to the node via more powerful nodes

This need for a flop-rich system led us to consider GPU-enabled nodes, due to their high flop rate performance—leading to the current Titan system

# 2013 Requirements Survey

## Topics:

- System hardware features needed
- Parallelism needs
- Use of heterogeneous compute nodes (e.g., GPUs)
- Characteristics of application code base
- Programming model
- Data storage, transfer and analysis issues

The community's ability to answer these questions as quantitatively as possible is extremely useful for our planning efforts

# Hardware features

How important is each of these hardware characteristics to your ability to deliver science:

1. **FLOPS** - floating point operations per second to perform calculations
2. **MEMORY CAPACITY** (more grid cells, particles, degrees of freedom., different algorithm, etc.)
3. **MEMORY BANDWIDTH** (sparse linear algebra, limited computations per accessed data element, etc.)
4. **MEMORY LATENCY** (unpredictable memory access patterns, unstructured grids, graph algorithms, etc.)
5. **INTERCONNECT BANDWIDTH** - the need to communicate large amounts of data between compute nodes
6. **INTERCONNECT LATENCY** - the need to communicate large numbers of messages between compute nodes
7. **LOCAL STORAGE CAPACITY** - increasing disk space for saved results and restart files
8. **DISK BANDWIDTH** - growing impact of data storage to disk on runtime
9. **DISK LATENCY** - large number of writes to disk of small size
10. **MEAN TIME TO INTERRUPT** - need for longer compute times between restart dumps
11. **ARCHIVAL STORAGE CAPACITY** - long-term storage of results
12. **WAN NETWORK BANDWIDTH** - need to communicate large amounts of data to/from offsite location

# Parallelism

- How much more available parallelism?
- How hard to extract?
- Do you rely on libraries for this?
- Which programming models (OpenMP, OpenACC, etc.) or which libraries (MAGMA, cuFFT, etc.) do you plan to use?
- Are there significant nonparallelizable parts of your code?

# Heterogeneous nodes

- Does your code run on heterogeneous hardware (NVIDIA GPUs, Intel Xeon Phi, etc.)?
- How hard is this port to do?

# Code base

- Do you expect to rewrite your code for new hardware?

# Programming model

- How adaptable is your code to new programming model?
- Do you anticipate:
  - Modifications of fundamental data structures, e.g., structure of arrays vs. array of structures?
  - Explicit memory hierarchy control, e.g., cache blocking?
  - CPU thread programming (OpenMP or Pthreads)?
  - Programming of accelerators with directives (OpenACC, Intel offload directives)?
  - Programming of accelerators with CUDA or OpenCL?
  - Explicit use of out-of-core techniques via disk or NVRAM?
  - More efficient checkpoint-restart and/or resiliency algorithms via utilization of NVRAM?
  - Explicit fault tolerance handling e.g., via FT-MPI to detect and correct failures?



# Programming model (cont.)

- Is performance portability to multiple platforms a challenge?
- How do you handle this?
- Do you use program development, optimization, and/or debugging tools today?
- What new capabilities do you anticipate needing for programming tools for using compute capabilities 10 times greater than available today?

# Data issues

- What are your data needs in 2016-2017 (scratch, archival)?
- Will you use in-situ data reduction, analysis or visualization?

# CORAL: the ORNL / LLNL / ANL Pre-Exascale Systems Procurement

- RFI at <http://www.csm.ornl.gov/CORAL-RFI/>
- Issues of concern mentioned in the RFI:
  - Resilience
  - Memory: improve memory bandwidth, mitigate memory capacity limitations
  - Programmability
  - Performance portability
  - Maximizing transfer rates between components of heterogeneous nodes
  - I/O system bandwidth to parallel file system
  - NVRAM approaches to persistent storage
  - Scalable networks
  - Embedded NICs
  - Advanced systems and power management features in the OS and runtime
  - Packaging (density, cooling, energy reduction, cable management)

# Questions?

Wayne Joubert

[joubert@ornl.gov](mailto:joubert@ornl.gov)

ORNL staff who provided input to this study:

Buddy Bland, Mark Fahey, Chris Fuson, Al Geist, Mitch Griffith, Rebecca Hartman-Baker, Joshua Hursey, Ricky Kendall, Don Maxwell, Bronson Messer, Maggie Miller, Hai Ah Nam, Jack Wells and Julia White

The research and activities described in this presentation were performed using the resources of the Oak Ridge Leadership Computing Facility at Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC0500OR22725.



# Supplementary slides

# OLCF Project Categories

OLCF requirements are driven by our projects which are subdivided into the following categories:

## 1. INCITE

- Large projects focusing on high-impact grand challenge research

## 2. ALCC

- High-risk, high-payoff simulations for special situations of interest

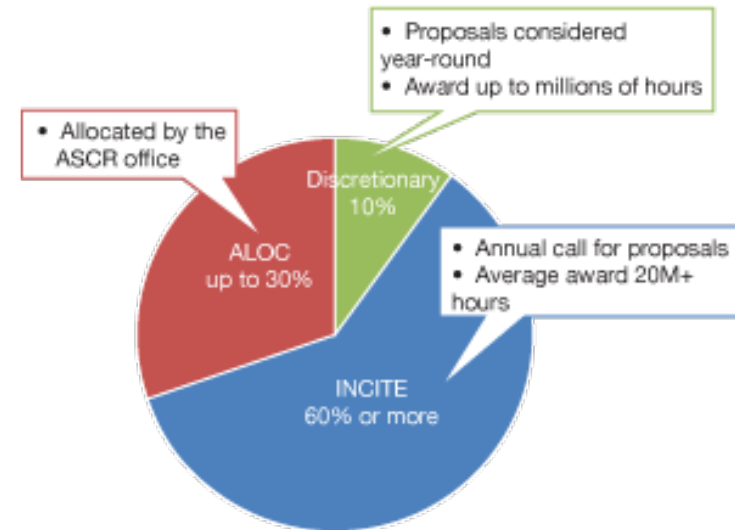
## 3. Director's Discretionary

- Small seed-time projects to jumpstart new efforts

The time-varying nature of the project mix from year to year makes precise requirements estimation a challenge—we do not have a fixed, recurring workload.



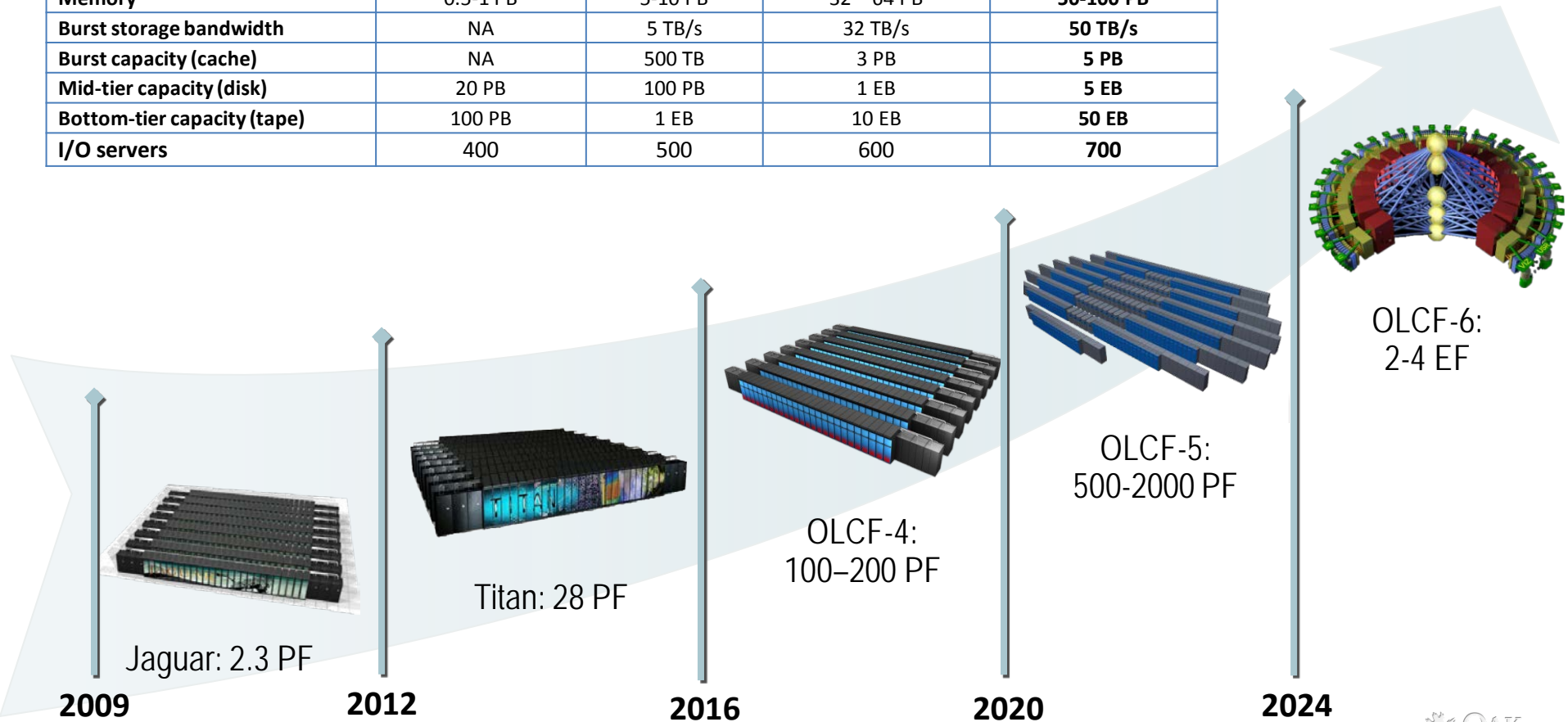
[www.doeleadershipcomputing.org/guide-to-hpc/](http://www.doeleadershipcomputing.org/guide-to-hpc/)



# The OLCF 10-Year Plan

- **OLCF has a 10-year plan to deploy and operate the computational resources required to tackle science problems of global importance**

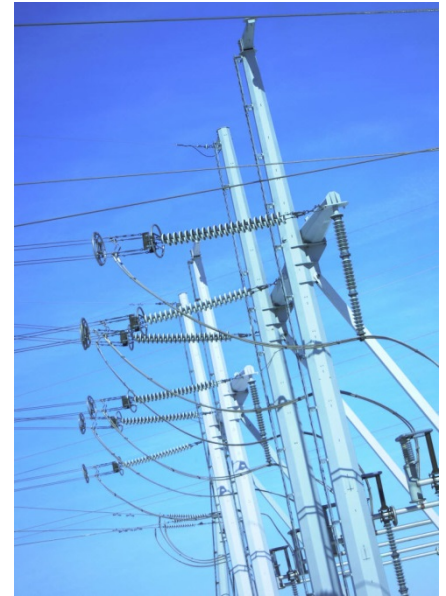
	2012	2016	2020	2024
<b>Peak flops</b>	10-20 PF	100-200 PF	500-2000 PF	<b>2000 - 4000 PF</b>
<b>Memory</b>	0.5-1 PB	5-10 PB	32 – 64 PB	<b>50-100 PB</b>
<b>Burst storage bandwidth</b>	NA	5 TB/s	32 TB/s	<b>50 TB/s</b>
<b>Burst capacity (cache)</b>	NA	500 TB	3 PB	<b>5 PB</b>
<b>Mid-tier capacity (disk)</b>	20 PB	100 PB	1 EB	<b>5 EB</b>
<b>Bottom-tier capacity (tape)</b>	100 PB	1 EB	10 EB	<b>50 EB</b>
<b>I/O servers</b>	400	500	600	<b>700</b>



# The Changing HPC Environment

## ➤ The power challenge

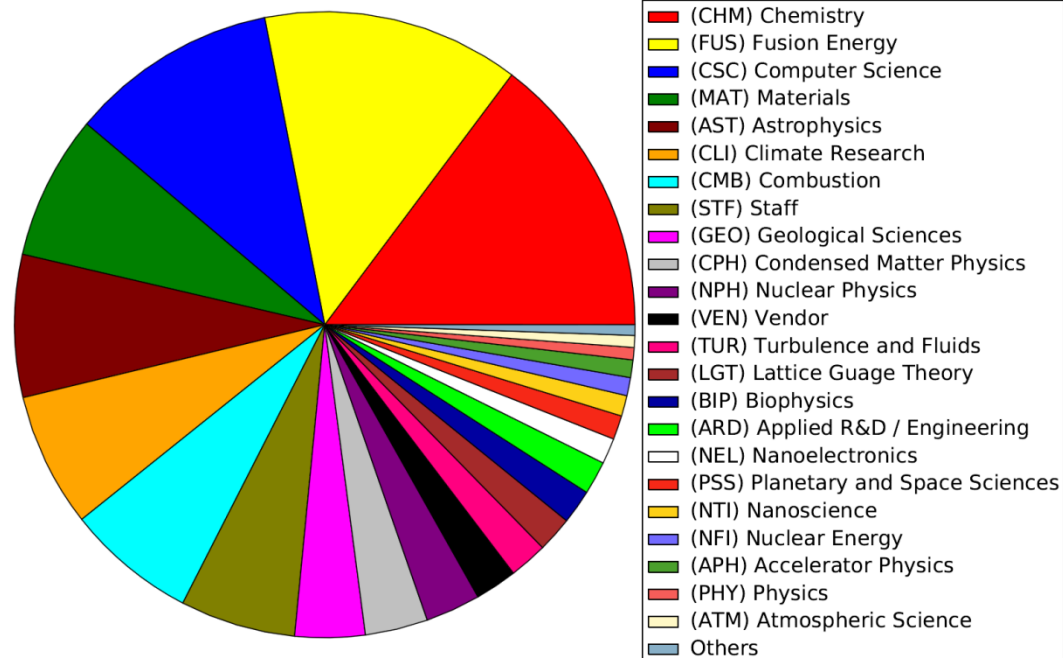
- OLCF and DOE are required to build systems within a plausible power envelope (~ 20-30 MW)
- Our goal is to deploy pre-exascale and exascale systems that address these challenges regarding hardware, software and power





# Workload of Past OLCF Systems

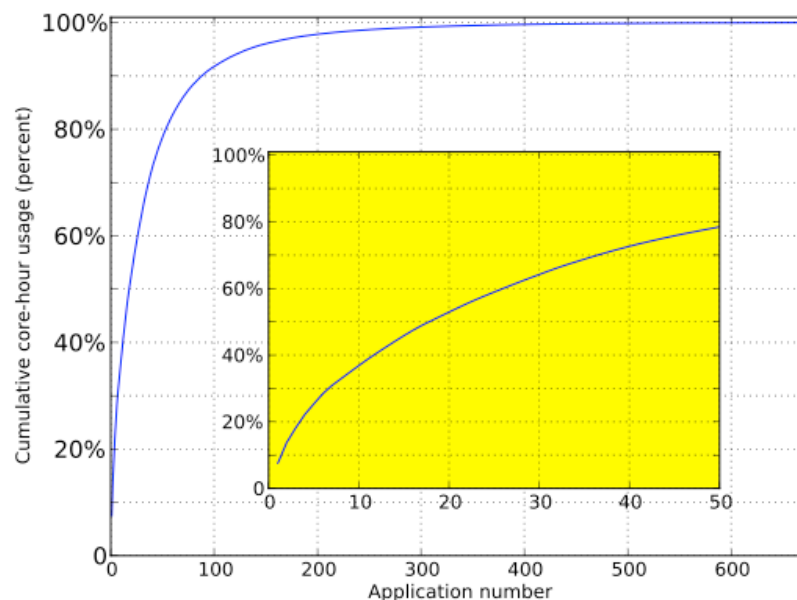
- Our systems must support diverse science domains with varied models and algorithms
- Requires a well-balanced system architecture



System core-hour usage by science domain  
JaguarPF, 2010-2011

# OLCF System Workload: Applications

- We tend to run a comparatively small number of codes that scale to large core counts on the system



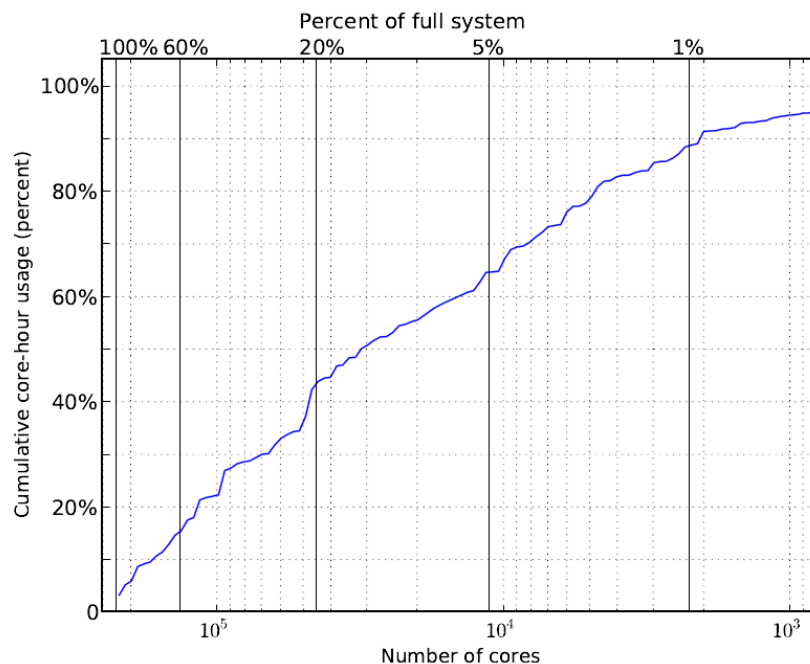
Cumulative system core-hour usage by application

JaguarPF, 2010-2011

The top 20 applications account for 50% of Jaguar usage

# OLCF System Workload: Scalability

- We attain a “leadership metric”—more than 40% of core-hours are spent in jobs run at 20% or more of full system size
- Users also run across the whole gamut of job sizes



Cumulative system core-hour usage by job size  
JaguarPF, 2010-2011

# OLCF System Workload: Algorithms

The most heavily-used codes are diverse in science domain and in algorithms used

Application	Primary Science Domain	Structured Grids	Unstructured Grids	FFT	Dense Linear Algebra	Sparse Linear Algebra	Particles	Monte Carlo
NWCHEM	Chemistry			X	X			
S3D	Combustion	X			X	X	X	
XGC	Fusion Energy		X				X	
CCSM	Climate Research	X		X		X		
CASINO	Condensed Matter Physics							X
VPIC	Fusion Energy	X					X	X
VASP	Materials			X	X			
MFDn	Nuclear Physics					X		
LSMS	Materials				X			X
GenASIS	Astrophysics		X			X		
MADNESS	Chemistry		X	X	X			
GTC	Fusion Energy	X				X	X	X
OMEN	Nanoelectronics	X				X		
Denovo	Nuclear Energy	X			X	X	X	X
CP2K	Chemistry	X				X	X	
CHIMERA	Astrophysics	X			X	X	X	
DCA++	Materials				X			X
LAMMPS	Chemistry	X		X			X	
DNS	Fluids and Turbulence	X			X	X	X	
PFLOTRAN	Geological Sciences	X	X		X	X		X
CAM	Climate Research	X		X	X	X	X	
QMCPACK	Materials						X	X
<b>TOTALS:</b>		<b>12</b>	<b>4</b>	<b>6</b>	<b>11</b>	<b>12</b>	<b>11</b>	<b>8</b>

Algorithm usage of top applications, JaguarPF, 2010-2011

# Data issues (cont.)

- How do you intend to share the scientific data emerging from your work in the 2016-2017 future?
  - Will you share your data with your scientific community?
  - Give an estimate of the useful lifetime of your scientific data?
  - What types of tools for data storage, movement, and analysis do you currently use?
  - Where do you see the need for tools development?
- What will be the total data set size for runtime checkpoint and restart?
- What will be the total data set size for post-runtime I/O?

# Recap

- Answers to these kinds of questions help us a great deal to understand your future system needs
- To the extent that you have this information or can determine it, it would provide useful input to us

# The OLCF Requirements Process

## Components:

- Science mission goals science drivers
- Requirements survey to elicit requirements from OLCF-supported projects
- Discussions with science code project leaders, team members and liaisons
- Information from OLCF project proposal applications
- Review of OLCF Leadership Computing usage logs over recent years to understand usage trends
- Survey of broader community and market trends
- Engagement with computer hardware and software vendors regarding capabilities of next-generation offerings and related trends