

The Hopper System: How the Largest* XE6 in the World Went From Requirements to Reality

Katie Antypas, Tina Butler, and Jonathan Carter CUG 2011, May 25th, 2011







Requirements to Reality

Develop RFP

Select vendor partner

Negotiate SOW

Deliver and Test System

Transition to Steady State







RFP Draws from User Requirements

- 13 'Minimum Requirements' (e.g., 24x7 support) that absolutely must be met
 - Proposals that don't meet are not responsive and are not evaluated further
- 38 'Performance Features' (e.g., fully featured development environment) wish list of features
 - Evaluated qualitatively via in-depth study of Offeror narrative.
- Benchmarks
- Supplier attributes (ability to produce/test, corporate risk, commitment to HPC, etc.)
- Cost of ownership (incl. life-cycle, facilities, base, and ongoing costs) and affordability
- Best Value Source Selection allows to evaluate and select the proposal that represents the *best value*







NERSC-6 Benchmarks









NERSC-6 SSP Metric

The largest concurrency time of each full application benchmark is used to calculate the SSP



For each benchmark measure

•FLOP counts on a reference system

•Wall clock run time on various systems





NERSC Cray Proposal is the Best Value

- Best application performance per dollar
- Highest sustained application performance commitment
- Best sustained application performance per MW
- Excellent in-house testing facility and benchmarking/performance/support expertise at Cray
- Easy to integrate into our facility
- Acceptable risk







Negotiation Challenges

- Cray proposed two technologies
 - XT5 available late 2009 with interconnect refresh 2010
 - Early cycles to users
 - Interconnect refresh incurs lengthy down time and hardware fallout
 - Older memory technology (DDR2)
 - Fewer cores per node
 - XE6 available mid 2010
 - Latest memory technology (DDR3)
 - Higher performance node
 - Latest interconnect delivered with the system
 - Delivered later





Feedback from NERSC Users was crucial to architecting Hopper

User Feedback

Hopper Enhancement

Login nodes need more memory

Connect NERSC Global FileSystem to compute nodes

Workflow models are limited by memory on MOM (host) nodes

8 external login nodes with 128 GB of memory (with swap space)

Global file system will be available to compute nodes

Increased # and amount of memory on MOM nodes
Phase 2 compute nodes can be repartitioned as MOM nodes







Feedback from NERSC users was crucial to architecting Hopper

User Feedback

Improve Stability and Reliability

Hopper Enhancement

•External login nodes will allow users to login, compile and submit jobs even when computational portion of the machine is down

•External file system will allow users to access files if the compute system is unavailable and will also give administrators more flexibility during system maintenances

•For Phase 2, Gemini interconnect has redundancy and adaptive routing.







Data and Batch Access





Hopper System

Phase 1 - XT5

- 668 nodes, 5,344 cores
- 2.4 GHz AMD 4-core Opteron
- 50 Tflop/s peak
- 5 Tflop/s SSP
- 11 TB DDR2 memory total
- Seastar2+ Interconnect
- 2 PB disk, 25 GB/s
- Air cooled

Phase 2 - XE6

- 6384 nodes, 153,600 cores
- 2.1 GHz AMD 12-core Opteron
- 1.27 Pflop/s peak
- 140 Tflop/s SSP
- 217 TB DDR3 memory total
- Gemini Interconnect
- 2 PB disk, 70 GB/s
- Liquid cooled









Hopper Phase 1 Installation





Delivery



Unwrap

Install



12



NERSC Site preparation



Unloading ...



Installation and Integration

from Tina Butler



Up and running!





Hopper places #5 in TOP 500 List at SC' 10







Hopper Early Hours

- ~320 million early hours delivered to science offices
- ~280 projects have used time
- ~1000 users have accessed the system

I.S. DEPARTMENT OF

JERG

 Consistently 300-400 unique users logged into system at any time

Office of

Science





Despite being a new, first-in-class peta-flop system, Hopper has run at a high utilization, with good stability from the start



•Over 81% utilization in the first month 2.5 month, (based on 24 hour day, including maintenances)

- •System problems that would have been full outages on the XT4 and XT5 can be contained on the XE6
- •Room for scheduling improvements, pack large jobs together, stabilize the system further
- •Maintenances a key source of lost utilization, look to minimize





Compared to the XT4 and XT5 most applications are seeing increased performance on Hopper

Hopper/Franklin and Hopper/Jaguar Performance Ratios



Below 1.0 - Application performs better on Hopper

•Applications run on Hopper, Franklin and Jaguar at same concurrency

•All benchmarks are pure MPI (except GAMESS which uses its own communication layer)

 Significant improvement on Hopper for GAMESS due to new Cray library on XE6 system. All other benchmarks use identical codes on Hopper. Jaguar and Franklin



Office of Science

Data from Nick Wright, Helen He and Marcus Wagner





Despite a slower clock speed, applications on Hopper perform better than the XT4 or XT5...

Metric	Franklin	Hopper	Impact on application performance
Proc clock speed	2.3 Ghz	2.1 Ghz	-
MPI latency	~6.5 us	1.6 us	•
MPI bandwidth	1.6 GB/sec/node	6.0 GB/sec/node	- -
Cache size	2 MB/socket shared L3	6 MB/6 cores shared L3	
Memory Speed	800 MHz	1333 MHz	•
Memory Bandwidth	~2 GB/sec/core	~2.2 GB/sec/core	\leftrightarrow

This is primarily due to the improved Gemini interconnect and thus less time spent in communication by applications







NERSC/Cray COE on Application Programming Models

GTC Fusion Application

pusher = shift = charge = poisson









Large Jobs are Running on the Hopper System

Breakdown of Computing Hours by Job Size



- Hopper is efficiently running jobs at all scales
- During availability period, over 50% of hours have been used for jobs larger than 16k cores.







Hopper is providing needed resources for DOE Scientists

- Over 320 M early hours delivered
- First time a peta-flop system is available to the general DOE research community
 - Production science runs
 - Code scalability testing
- Hopper is a resilient system
- Component failures are more easily isolated
- Survives problems that case full crashes on XT4 and XT5
- Researchers appreciate the stability of the system and they want more time



"The best part of Hopper is the ability to put previously unavailable computing resources towards investigations that would otherwise be unapproachable." – Hopper User





Acknowledgements

- This work was supported by the Director, Office of Science, Office of Advanced Scientific Computing Research, of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231.
- The authors would like to thank: Nick Wright and Helen He for providing early COE results; Manuel Vigil and the ACES team for valuable discussions and test time at the factory; and the Cray on-site and Cray Custom Engineering staff for valuable discussions.

