



**CLUSTER**<sup>®</sup>  
**RESOURCES** INC.

## **Moab and TORQUE Achieve High Utilization of Flagship NERSC XT4 System**

Michael Jackson, President  
Cluster Resources  
[michael@clusterresources.com](mailto:michael@clusterresources.com)  
+1 (801) 717-3722

# Contents

1. Introduction
2. Motivation
3. Design Principles
4. Required Attributes
5. Conclusion



# NERSC

- While NERSC provides some of the largest computing and storage systems available anywhere, what distinguishes NERSC is its success in creating an environment that makes these resources effective for scientific research.
- US Department of Energy
- 19,000+ Cores
- 3,000+ Computational Scientists

# Cluster Resources Technology Focus

## Technology Areas

- Clusters
- Supercomputers
- Grids – Local Area Grid (One Administrative Domain)
- Grids – Wide Area Grid (Multiple Administrative, Data and Security Domains)
- Data Centers
- Hosting
- Utility/Adaptive Computing

## Platforms

- XT3 / XT4 ...
- X1E
- XD1
- Other Platforms



# Moab

## Integrates

- Scheduling
- Monitoring from Multiple Sources
- Managing Multiple Resources and Environment Factors
  - Richest Set of Policy Options
  - Highest Utilization
  - Better Behaviour
- Reporting on Usage, Behaviour, etc.
- Events Inside and Outside of the Scheduler

# TORQUE

## Aspects

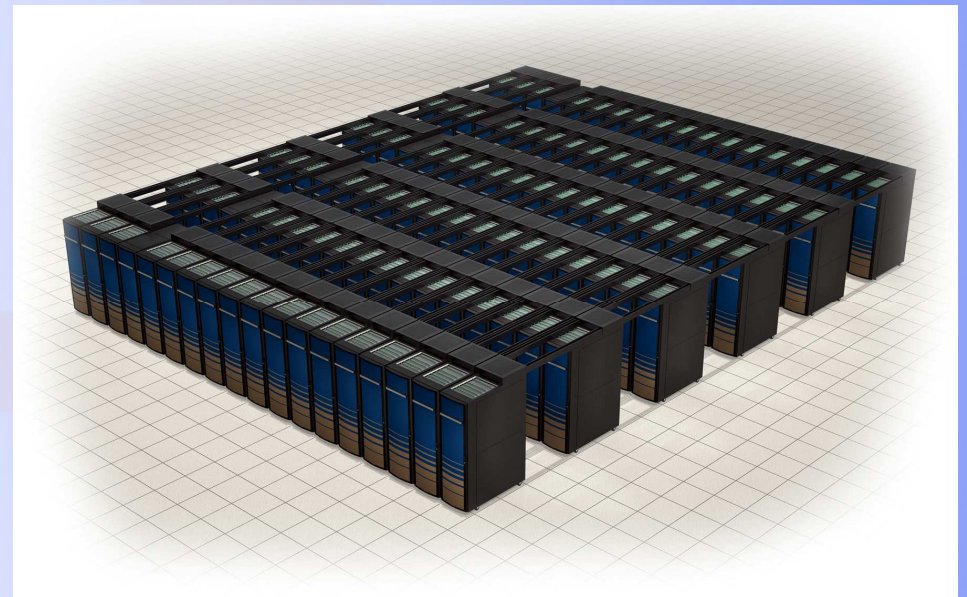
- Used by many thousands of sites
- Added High Availability (Active Active Active)
- Added Job Array Support
- Monitor
- Queue
- Execute Jobs

## The Motivation

- Interoperability with Cray's CNL (a.k.a. Cray Linux Environment) and ALPS
- Ultra High Efficiency (90 – 99%) with common achievement of 97+%
- Dynamically Modifiable Limits, Policies and Management
- Intelligent and Hierarchically Inheritable Policies
- Advance Reservations
- High Availability
- Holistic Reporting and Integration with Custom and Open Source Accounting Tools
- Support for Heterogeneous Resources
- Grid Support
- Rich Configuration Options
  - Dynamic Backfill
  - Fairshare
  - Preemption
  - Rich Credentials

## Design Principles (1 of 4)

- Highly Flexible
- Easy to Manage
- Low User Knowledge Requirement
- Efficient and Effective in their use of the System
  
- Credentials
  - Standard Users
    - Limit Excessive Usage
  - Special Users
    - Special Rights
  - Staff Members
    - Added Rights
  - Administrators
    - Extensive Control



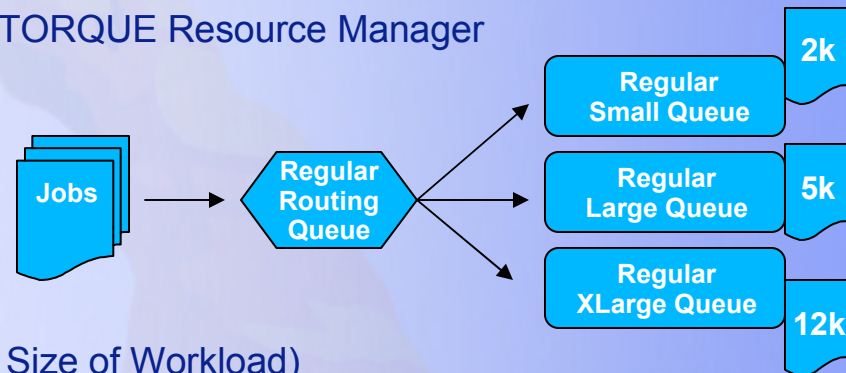


## Design Principles (2 of 4)

- Queue / Class

- Set up Queues and Routing Queues in TORQUE Resource Manager

- Regular Routing Queue
    - Special Routing Queue
    - Interactive Queue
    - Debug Queue
    - Admin Queue
    - Transfer Queue (For data transfer)



- Set up Sub Queues in Moab (Based on Size of Workload)

- Regular-Small (as large as 2,000 cores)
    - Regular-Large
    - Regular-Extra\_Large
    - Etc.
      - No added work by user to specify sub queue, job is allocated to proper sub-queue automatically
      - Large jobs were favored (so small jobs did not push them out of the way)
      - Small jobs already get the benefit of the dynamic back fill
  - Transfer Queue
    - 0 Tasks (No processing, only allow data transfer and other such services)
    - Very high priority (get movement accomplished now)
    - Short duration allowed
    - Tracked for accounting purposes

## Design Principles (3 of 4)

- Quality of Service (QoS)
  - Apply rules/policies that can be inherited by multiple entities easily
  - Avoid multiplication effect on queues/classes
- Prioritization
  - NERSC implemented a design that focused on Simplicity
    - 1 = 1 minute
    - 60 = 1 hour
    - 1440 = 1 day (60 minutes times 24 hours)
    - NERSC selected 1440 as the default value and other modifiers applied against this base
  - Example
    - 0 for low priority
    - 2 for regular priority
    - 5 for high priority
  - If a low, regular and high priority job were run on the same day, the high priority job would look like it had been in the queue for 5 days, the regular for 2 and the low as just submitted. Therefore if a high priority job is submitted two days after a regular priority job, it still will be favoured to run first.
  - A high priority job would run in front of a regular job that had been in the queue for less than three days ( $5 - 2 = 3$ ) and a low priority job that had been in the queue for less than five days ( $5 - 0 = 5$ ). Small jobs already get the benefit of the dynamic back fill

## Design Principles (4 of 4)

- Accounts and Groups
  - Set in place to make quick adjustments
  - Allows for accounting in the meantime
- Periodic Standing Reservations
  - Business Days and Business Hours Availability
    - Interactive Sessions
    - Data Transfer
    - Debugging
  - Optimization of Resources (Dissolve the reservation after hours to resources are used after hours)
  - Examples
    - Monday through Friday each week
    - 8:00 AM to 8:00 PM
  
    - Tuesday and Thursday
    - 2:00 PM to 4:00 PM

## Required Attributes (1 of 4)

- Ultra High Node Utilization
  - Utilization
  - Packing
  - Targeted Usage to What is Important
  - Reports for Evidence
  
- NERSC
  - Mix of Workload
  - 19,000 + (single system jobs)
  - Avoid fragmentation caused by some resource managers
  - Use policies not silos
  
- Maintenance Reservations
  - Don't "drain the system", nor send out emails to stop working or start working
    - Long periods
    - Wasted resources
  - Maintenance reservations run all workload that can complete in time
  - No notifications required
  - Partial maintenance reservations also allowed (Moab applies work to all other areas)

## Required Attributes (2 of 4)

- Ultra High Node Utilization
  - NERSC is commonly maintaining 97+% node utilization
  - Checkpoint Restart
    - Improves reliability
    - Improves efficiency when used in conjunction with maintenance reservations
    - Improves ability to protect and manage long running jobs as well as reduce the work of breaking large running jobs into smaller jobs
  - Dynamic Backfill (One example of packing)
    - Dynamic means it adapts each iteration to optimize within what is now available
      - Jobs finish early
      - Jobs get cancelled
      - Etc.
    - First Fit – (NERSC's choice) Focuses on fairness and prioritization
    - Best Fit – Applies largest single job to available space
    - Greedy – Applies largest combination of jobs to fit available space
  - Dynamic backfill is increasingly important the larger your jobs are
  - Other Methods of Packing

## Required Attributes (3 of 4)

- High Availability
  - Moab high availability
  - TORQUE Resource Manager
    - Active-Active-Active High Availability
  - Checkpoint restart
  - Intelligent resource allocation (Work around full or partial failures)
    - Lustre clean up cycles
    - Failed nodes
    - Etc.
- Other Attributes
  - Highly flexible management (High levels of configuration, nearly all dynamic)
  - Ease of use and ease of administration
    - Inheritable rules – no user work
    - Routing queues – no user work
    - Web based job portal
    - Workload templates
    - Graphical administration/helpdesk
    - Reporting tools

## Required Attributes (4 of 4)

- Accounting and Reporting
  - NERSC Information Manager (NIM) maintains what the users or projects have left to use
  - TORQUE Resource Manager provides usage tracking details
  - Each day the system evaluates what was used and debits against what is left
  - If the user is out of credits, no new workload can be submitted
  - Moab is able to provide rich visual charts and reports
  - Comparable to Gold Allocation Manager ([www.clusterresources.com](http://www.clusterresources.com))
    - Real-time tracking
    - Block or lower priority
    - Account across complex set of users, projects, other sites, etc.

## Conclusion

- NERSC effectively applied advanced scheduling and workload management
  - Highest Utilization
  - Higher Flexibility
  - More Usable System
  - Lowered Administrative Requirement
  
- Made possible from Moab & TORQUE capabilities, but ultimately resulted from the effective design choices of NERSC
  
- Acknowledgments
  - Nicholas P. Cardo – XT4 System Lead – NERSC – Lawrence Berkeley National Laboratory
  - Scott Jackson – Core Developer for Cray Platforms – Cluster Resources, Inc.
  - Craig Laurence – Technical Editor – Cluster Resources, Inc.
  
- Send Questions to Author:
  - Michael A. Jackson: +1 801 717 3722, [michael@clusterresources.com](mailto:michael@clusterresources.com)



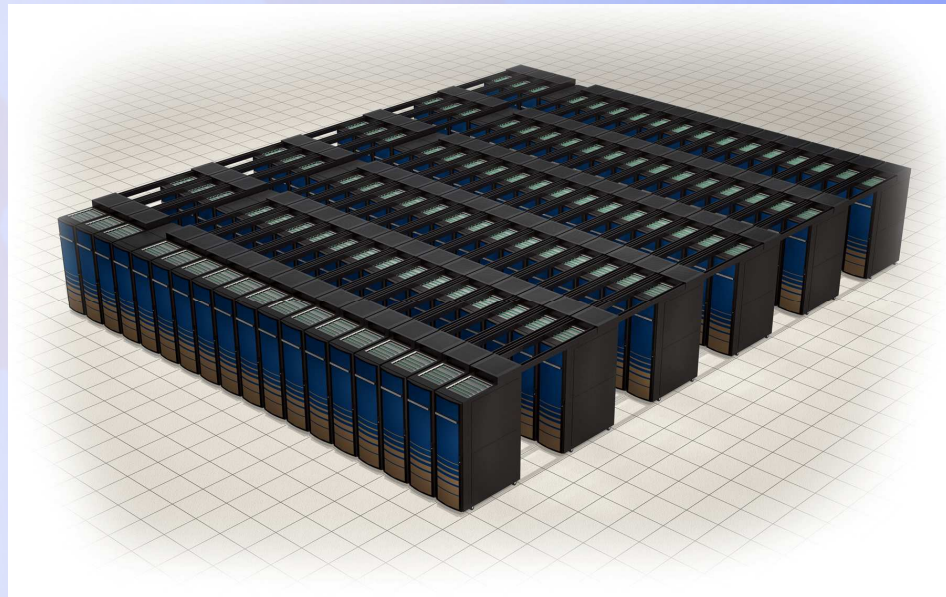
# Appendix

# Managing Leadership Systems w/ Moab

## NERSC

**Over 19,000 cores**  
**Cray XT4**

- AMD Opteron™
- ~100 racks



## Awarded Largest HPC Management Contract in History

# US Department of Energy

250,000+ processors, 300+ clusters

“Partnerships such as this one are a key element of the ASC Program’s success in pushing the frontiers of high performance scientific computing. Only by working with leading innovators in HPC can we develop and maintain the large scale systems and increasingly complex simulation environments vital to our national security missions.”

—Lawrence Livermore National Laboratory



## Managing Leadership Systems w/ Moab

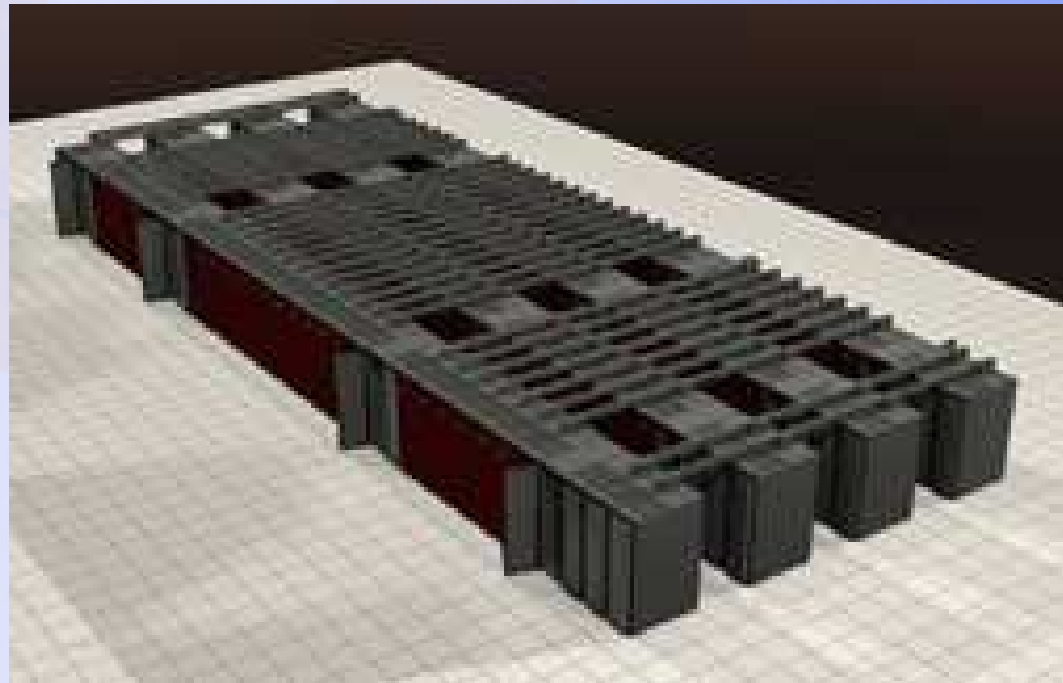
# Sandia – Red Storm

### Red Storm:

12,960 CPUs

**Cray XT3**

- 124.42 teraOPS theoretical peak performance
- 135 racks
- AMD Opteron™
- 40 terabytes of DDR memory
- 340 terabytes of disk storage
- Linux/Catamount OS
- <2.5 megawatts power & cooling



# Managing Leadership Systems w/ Moab

## ORNL

**Jaguar**: ~18,000  
core **Cray XT3**  
moving to  
1 Petaflop

**Phoenix**: 1,024  
core **Cray X1E**

**RAM**: 256 CPU  
**SGI Altix**



## Unified Management & Optimization

# Barcelona Supercomputing Center



## Europe's Largest Supercomputer/Cluster

5<sup>th</sup> Largest HPC System in the World

## Top500 Systems licensed with Moab Leadership Class Systems

**#1, 2, 4, 5, 6, 10,  
11, 19, 20, 23, 24, 25,  
28, 35, 37, 44, 51, 54, 58,  
60, 62, 69, 71, 75, 78, 79, 84, 87,  
96, 114, 140, 142, 143, 146, 176, 202, 203,  
230, 296, 297, 356, 366, 412, 445, 450, 451, 454, 488**

# Company: Moab in the HPC Stack

