

## HPC Annual Report 2006-2007

Aug 16, 2007

### Overview

A lot of information about the University of Florida High Performance Computing Center and its activities is available on the web at <http://www.hpc.ufl.edu>. This report provides a brief overview of the activities of the Center and its governing body the HPC Committee (ITAC-HPC). Because this is the first annual report, some information from previous years is included for completeness.

First, the state of the HPC Center is described, including personnel, hardware, software, services, user community and an overview of the budget. Then we list some of the accomplishments and activities of the last year. Because this is the first year the report is published, we include some information from earlier years. Next, we present the plans for the 2007-2008 year. Finally, we show some graphs with usage statistics of the cluster over the year.

### State of the HPC Center

**Phase I** The HPC Center started operation, after several years of preparation by the HPC Committee, with Phase I of the HPC cluster which became operational in August of 2004. The cluster was manufactured by Dell and had 200 compute nodes with 2 GB of RAM and two Xeon processors in each node. Four server nodes provided access to 32 TB of storage and two tape robots for backup. All nodes were connected by a CISCO 6509 Gigabit Ethernet switch. The system was installed in the New Physics Building machine room 2250. The cooling capacity of the room had to be upgraded for it to be able to cool the new cluster and this caused a six-month delay of the start of operations. The cluster was being taken into use by UF researchers very quickly and ran at capacity until it was taken out of service to make room for the Phase II cluster in September 2005. The job scheduler was PBSPro which provides very simple scheduling that satisfied the need at the time.

**Phase IIa** In August 2004, NSF awarded a proposal requesting funding for building a high speed research network on campus and for research in campus-wide accessible storage made available over that network. The principal investigators are S. Ranka, P. Avery, A. George, S. Trickey, P. Sheng. A Collaborative Research Agreement with CISCO was reached during the summer of 2005. This agreement included the purchase of Gigabit Ethernet switches from CISCO to create a 20 Gigabit/sec campus research network connecting several buildings, including the machine room holding the HPC cluster, to SSRB and the Florida Lambda Rail and the National Lambda Rail and UltraLight. The agreement also included a significant donation from CISCO of InfiniBand equipment to provide the fast interconnect for the Phase II cluster. The Phase II cluster fits in the same footprint (space, electrical, and heat) of the Phase I cluster. It was manufactured by Rackable and has 200 nodes with 4 GB RAM and four Opteron CPUs on two sockets. All nodes are connected by InfiniBand. Six servers provide access

to 40 TB of storage. The installation was started in Dec 2005 and the system became fully operational in May 2006. During this six-month period, it was possible for “friendly users” to use the system, but there were some system issues that needed to be resolved. One issue was to build a fast, parallel file system for the nodes. In the end the RapidScale file system from Rackable was deployed.<sup>1</sup> The new cluster, four times as powerful as the Dell cluster it replaced, quickly filled up to utilization of over 95% most of the time. The PBSPros scheduler was also being used on the Phase II cluster.

**Phase IIb** In December 2006, an expansion of the Phase II cluster was configured and ordered. This expansion doubles the number of nodes to 400, bringing the total number of CPUs to 1,600. The nodes of the expansion have the same CPUs but have 8 GB of RAM. Two thirds of the added nodes will be joined into the InfiniBand fabric, with 80 nodes configured for Gigabit-Ethernet-only connectivity to primarily serve serial jobs. To house the expansion, the Tier 2 cluster was moved into the QTP machine room NPB 1114 to make room for two racks in NPB 2250; the third, Ethernet-only rack was also placed in NPB 1114. The installation of the additional nodes did not require any system downtime and the expanded cluster became available to users in February 2007.

**Scheduler** In April, a new scheduler was installed and tested that allows more flexibility to meet the requirements for quality of service to investors as specified in the HPC Contract documents. The new scheduler is called the Maui scheduler. It is being used at many HPC centers and developed and maintained by a large group of system administrators. The new scheduler works much better, but our testing and exploring of the capabilities is not yet complete. At this time, the scheduler does not do exactly what the documents specify, but it does come close to the meaning of the specification. The HPC Center is still able to provide, in a timely fashion, to all researchers on campus more cycles than they need.<sup>2</sup>

**Uptime** There was a period of scheduled downtime in September 2006, one during the second week of January 2007, and one during the first week of August 2007. These periods are used to install software updates and make small hardware repairs and changes as needed. Between the last two periods of scheduled downtime, there were two or three occasions where a problem with hardware or software caused the cluster to be unavailable for a short time or created a condition that required that every node in the cluster be rebooted. Otherwise the cluster was operational 24 hours per day, 7 days per week.

**Research Network** The cluster became a part of the Open Science Grid (OSG) and is being used by OSG users. The connectivity via the Campus Research Network to the Florida and National Lambda Rail is actively being used by the High Energy Physics Tier 2 group, and by Prof. Haselbacher (MAE) to move large data sets between the cluster and the storage robot at NCSA.

---

<sup>1</sup> See the “Parallel file system” item in the section “Accomplishments and activities” to find out about TerraScale and the Terragrid file system and the acquisition by Rackable.

<sup>2</sup> This, of course, is not necessarily as many cycles as the researchers could use.

## Budget overview

The HPC Center operates on a very small annual budget for operations, but it manages a significant amount of assets and funds. A brief overview is presented in the table below.

Description	2003-2004	2004-2005	2005-2006	2006-2007
Faculty	\$200,000	\$1,092,000	\$0	\$620,000
Department	\$0	\$213,000	\$25,000	\$67,000
College	\$100,000	\$463,000	\$130,000	\$182,000
OIT	\$370,000	\$301,000	\$110,000	\$336,000
<b>TOTAL IN</b>	<b>\$670,000</b>	<b>\$2,069,000</b>	<b>\$265,000</b>	<b>\$1,205,000</b>
Staff	\$0	\$171,000	\$244,000	\$263,000
Facilities	\$70,000	\$0	\$0	\$158,000
Equipment	\$600,000	\$1,667,000	\$0	\$760,000
Operation	\$0	\$10,000	\$10,000	\$61,000
<b>TOTAL OUT</b>	<b>\$670,000</b>	<b>\$1,848,000</b>	<b>\$254,000</b>	<b>\$1,242,000</b>

## List of investors

Some investors are investing as head of a college or department, some as individual faculty, some are playing both roles.

Asthaagiri A., professor, Chemical Engineering

Avery P., professor, Physics

Balachandar S., chair, Mechanical and Aerospace Engineering

Cheng H.-P., professor, Physics

Curtis J., chair, Chemical Engineering

Fortes J., Electrical and Computer Engineering

Glover J., dean and Sabin J., associate dean, College of Liberal Arts and Sciences

Haselbacher A., professor, Mechanical and Aerospace Engineering

Hoit M., CIO, Office of Information Technology

Khargonekar P., dean, College of Engineering

Law M., chair, Electrical and Computer Engineering

Philpott S., professor, Materials Science and Engineering

Shea J., professor, Electrical and Computer Engineering

Sinnott S., professor, Materials Science and Engineering

## Staff

The center has three permanent staff members:

1. Charles Taylor, associate director, TEAMS position funded by OIT
2. Craig Prescott, Scientist position funded by CLAS, being converted to TEAMS
3. Jon Akers, TEAMS position funded by ECE

In addition, the center has had several students working part-time. The funding varies. CLAS supported a student for some time, OSG/Tier2 has supported several students,

## Annual Report 06-07

including one now. One student completed her Master's degree in engineering in August 2006.

### **Accomplishments and activities**

**Parallel file system** One of the goals of the HPC effort at UF is to create uniform and high performance access to data storage. The first component is to create access from all nodes to a uniform file system. To get high performance data throughput requires a parallel file system. HPC Center engineers investigated several such file systems, e.g. Lustre and Panasas, and tried some on the Phase I cluster, e.g. PFS. For the Phase II cluster, two vendors were invited to install software and hardware for evaluation during winter and spring of 2006: CrossWalk and TerraScale. The HPC Center staff worked with the engineers of both companies to try and reach performance that was reasonably close to the maximum achievable with the available hardware. The engineers from TerraScale addressed several problems in the TerraGrid parallel file system. Ultimately, excellent input and output performance was achieved with this system. The HPC Center decided to install the TerraGrid file system in April 2006. As a result of the outstanding performance results achieved at UF, Rackable decided to acquire TerraScale in the summer of 2006. The file system is now known as the RapidScale file system from Rackable.

**CISCO Collaborative Research Meetings** As part of the 2004 NSF MRI award to build a high performance storage infrastructure on the UF campus, CISCO donated a significant amount of hardware for the Phase II cluster. The donation was part of a “Collaborative Research Agreement” between CISCO and UF. As part of the agreement two meetings were held at UF, the first on Feb 9 and 10, 2006 and the second on Sep 12 and 13, 2006. The agenda of the first meeting was dominated by technical issues related to the installation and performance of the, then new, Phase II cluster and its InfiniBand fabric. The architecture of the cluster and of its storage system and of the 20 Gb per sec Research Network and of the campus grid were also discussed in detail. The second meeting focused on the actual research that was being done with the HPC cluster, with its InfiniBand fabric supporting parallel applications and access to the parallel files system for high performance input and output, and with the UF Research Network, which is powered by the CISCO Ethernet switches. The lessons learned from the installation of the cluster and of the InfiniBand fabric were also discussed in detail. Detailed plans were made for joint activities and publicity by CISCO and UF at Super Computing 2006.

**Super Computing 2006** The second component of the HPC Center goal regarding data storage is to create a file system that can be accessed from all clusters on the campus grid that are connected by the 20 Gb/s Research Network. During the research for technology that can support this, Charlie Taylor and Craig Prescott found and contacted a small company in Canada called Obsidian Research. This company makes a long-distance extender for InfiniBand fabrics. A test unit was sent to UF in August for a research project with results to be shown at the annual Super Computing conference scheduled in Tampa in November 2006. A dedicated fiber was activated between the HPC Center machine room in NPB and Alan George’s HCS laboratory in Larsen Hall. Nodes in the HCS cluster were configured as InfiniBand clients of the parallel file system exported by the HPC cluster. Performance was measured in a series of experiments. The experiments were successful and it was decided to try and make a demo where a computer in the

exhibition hall in Tampa at Super Computing 2006 would access the servers in Gainesville over the Florida Lambda Rail. With help from many people in many organizations (CISCO, UF, Obsidian, Rackable), including Chris Griffin and Dave Pokorney from CNS, this feat was accomplished.

**Obsidian Research paper** The work done on the UF campus with the LongBow equipment from Obsidian Research was written up in detail in a paper “Comparative Performance Analysis of Obsidian Longbow InfiniBand Range-Extension Technology”, by Craig Prescott and Charles Taylor to be submitted to the IEEE Computer journal.

**Storage World 2007** The University of Florida was one of the five finalists in the category of “Innovation and Promise” in the Storage World “Best Practices in Storage” Awards Program. The award identifies and acknowledges excellence among users of storage IT solutions and approaches. Finalists in each category were honored in a ceremony April 18, 2007 at the Storage Networking World conference in San Diego. All five finalists received an award. Jon Akers attended to receive the award, with travel expenses paid by Rackable and registration paid by StorageWorld.

**HPC Contract** During the academic year 2006-2007, the ITAC-HPC committee focused its monthly meetings on drafting, discussing and modifying a set of documents that describe the process that should make the HPC effort at UF sustainable. The result is a set of three documents: HPC Contract Summary, HPC Sustainability Plan, HPC Contract. The first document is a one-page summary, the second provides the definitions and arguments behind the contract, the third is the actual contract between the stake holders of HPC at UF. The stake holders are defined to be the faculty and their research associates, the administration (DDD), and the CIO and the HPC Center. The documents were adopted in April 2007. In addition the committee adopted the HPC Center Management Plan.

## Projects and plans

**Campus storage server** Using the storage part of the 2004 NSF MRI award that funded the creation of the Campus Research Network, the HPC Center plans to buy and install a second storage server for the campus grid, the first server being the HPC cluster. The latest release of the RapidScale files system, being installed during the August 2007 maintenance window, allows both InfiniBand and Ethernet clients to access the same file system at the same time. Thus clusters will be able to mount the HPC cluster file system across the 20 Gbps Research Network. The execution of this plan has been delayed because of accounting problems.

**Small-Tree contract** Charlie Taylor has successfully negotiated a multi-year contract with a software company, called Small-Tree. The contract will pay for several graduate students to work in the HPC Center. The company focuses on software for Macs and the HPC part of the contract is to develop an open source version of the software for the Linux system.

**Scheduler** The demands on scheduling the 1600 CPUs to meet the complex and varying requests of the UF research community are very high. The new scheduler does much better than the old PBSPro scheduler, but a lot of work needs to be done still to allow predictable start times for parallel jobs.

**Campus grid** We plan to bring other clusters on campus into the grid structure with the HPC cluster and the Tier 2 cluster. This poses many challenges, both of a practical and of a political nature.

**Preparing for phase III** An awareness campaign in the Health Science Center will be launched with the goal to identify the needs of the HSC for HPC.

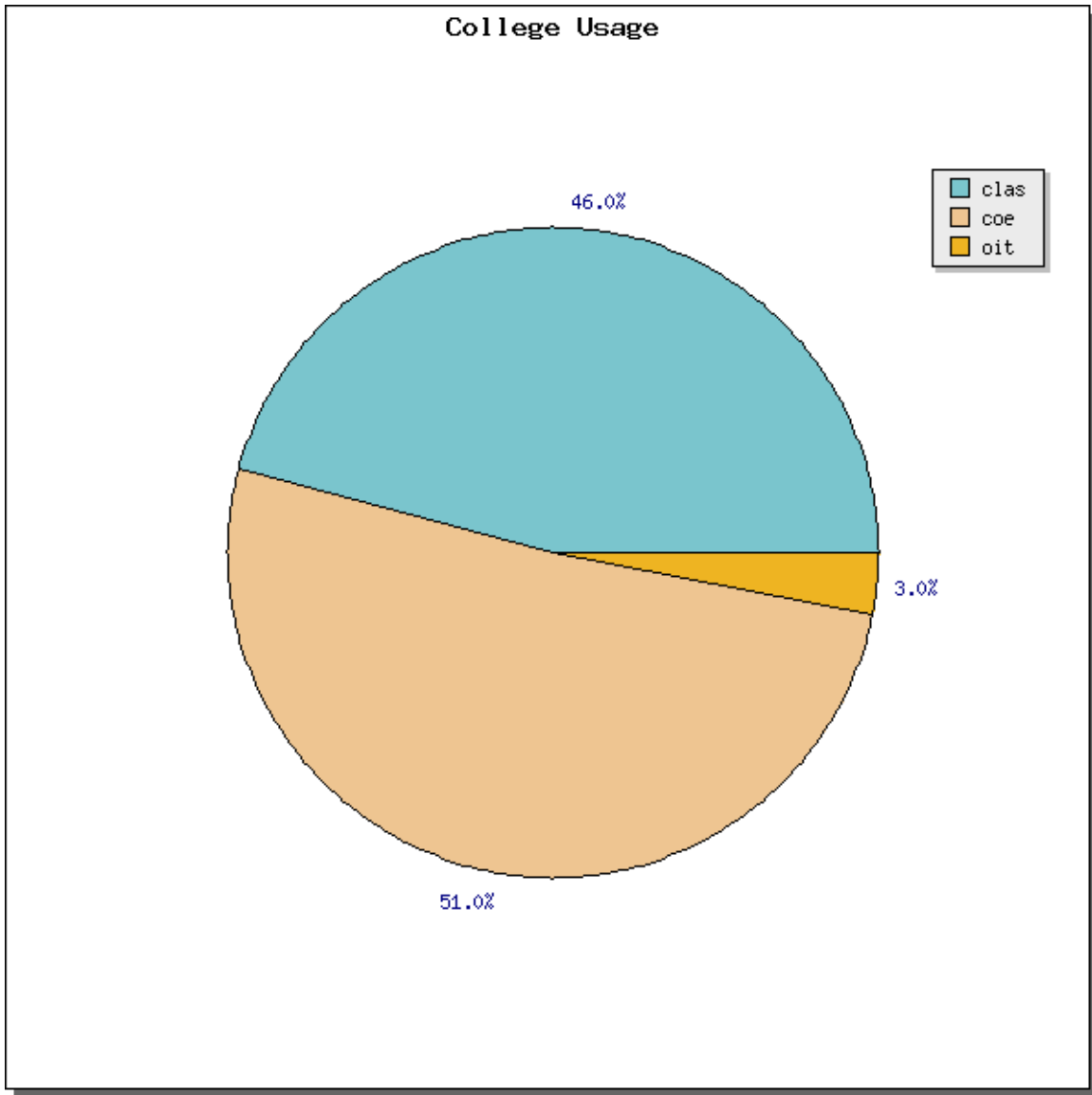
### Usage statistics

The HPC Center website contains a lot of information, including live information about the utilization of the system: <http://www.hpc.ufl.edu/index.php?body=util>. There are several displays worth pointing out.

1. “CPU Usage Summary“ <http://www.hpc.ufl.edu/index.php?body=pbs/nodestate> shows in a single view each of the 1600 CPUs with a color code indicating whether the CPU is idle or busy on a serial job, or a parallel job with 2-8, 9-32, 33-128, or 129 or more processors.
2. “CPU Job Utilization” <http://www.hpc.ufl.edu/index.php?body=pbs/cpustat> summarizes the same information with a bar chart and a pie chart.
3. “Torque Queue Status” and “Maui Queue Status” show a full list of all jobs waiting and running with details such as number of CPUs requested and time accumulated in the queue or time accumulated executing.
4. “Cluster Usage Statistics (last 7 days)” show the percentage of the cluster time used labeled by research group as a pie chart and as a table. You can also request this information for a different number of days than 7.
5. “Usage by College” and “Usage by Department” show the percentage labeled by College and Department for the last 30 days. You can request any number of days at the bottom of the screen. The plots for 365 days before Aug 14, 2007 are included below.

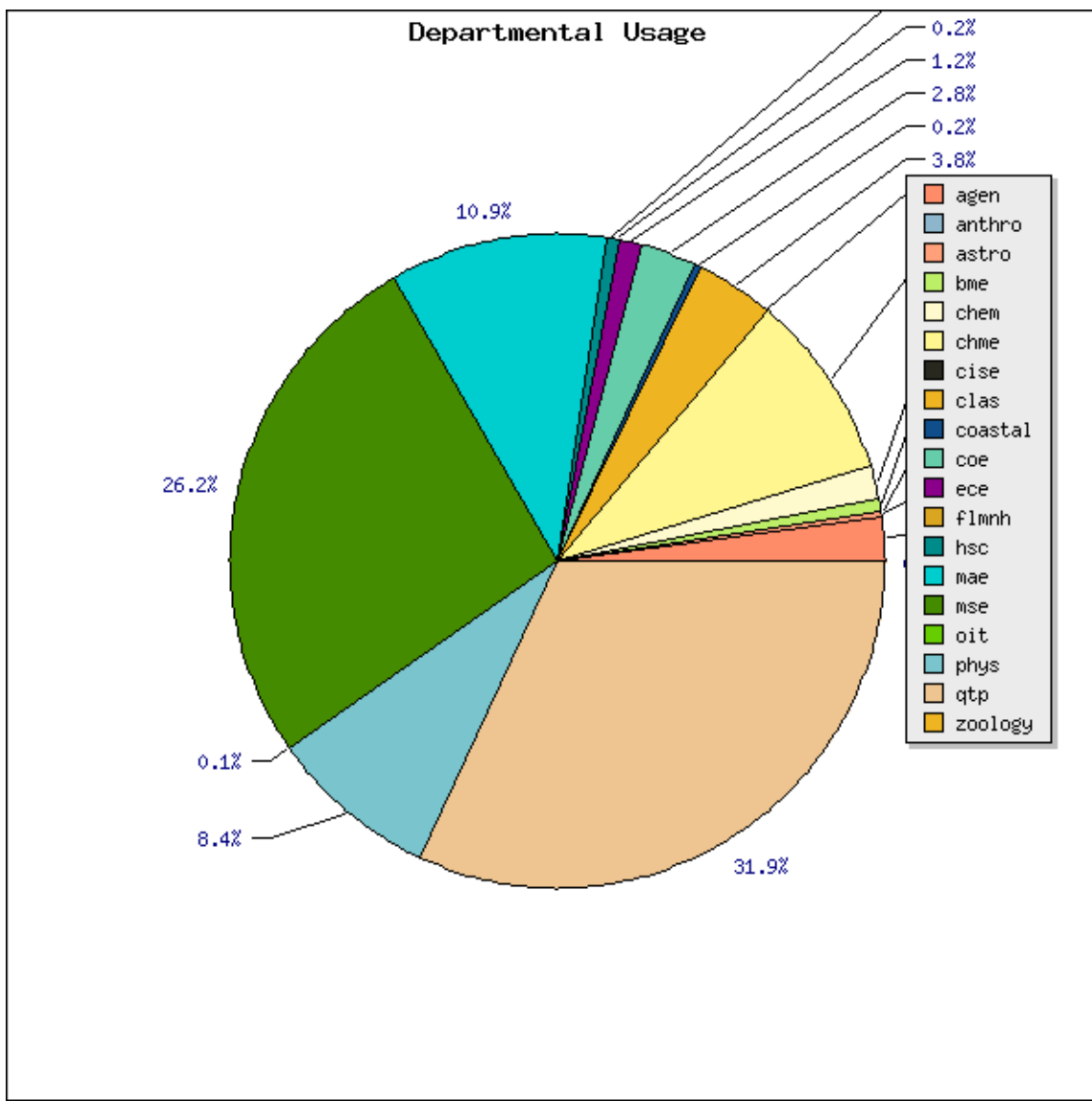


Usage by College from Aug 15, 2006 through Aug 14, 2007



The College named "OIT" encompasses all users on campus who do not belong to CLAS or COE, which are the Colleges that have invested in the HPC Center. All these users run under the OIT investment priority.

Usage by Department from Aug 15, 2006 through Aug 15, 2007



The Department named “OIT” encompasses all users on campus who do not belong to a Department or College that has invested in the HPC Center.

## **Work done during the maintenance window**

During the period July 31 through Aug 9, the HPC cluster was down for maintenance work, installation and reconfiguration of hardware and software.

### **Hardware changes**

- Added memory to imgsrv
- Installed IB card in imgsrv
- Changed imgsrv CPUs from 2x248s to 2x265s
- Replaced RAID card in racksrv
- Replaced power supply in racksrv
- Replaced motherboard in racksrv
- Replaced memory on hpcio5
- Upgraded submit from 4-way ASUS K8N-DRE to 8-way Tyan m4881
- Replaced memory on submit
- Replaced power supply on hpcio2
- Install 10-gigabit ethernet card in submit
- Rerouted I/O Node IB Cables
- Recabled IPoIB bridge-groups to corresponding 6506 port-channels
- Installed 6704 blade into the 6506
- Wired the 4948 to the 6506 via two 10-gbit links

### **Network changes**

- Attached Cisco 4948 to 6506 via ISL
- Configured NAT on 6506
- Configured ACLs on 6506
- Configured port-channels for bond0 interfaces on hpcio1 - hpcio8
- Created port-channel for bond0 interface on altix to ethsw4948
- Placed ib1 interface of each I/O node on separate subnets for IPoIB
- Configured bridge-groups 41-48 on IPoIB gateway
- Configured port-channels 41-48 for corresponding IPoIB gateway bridge groups
- Moved ISL to Netgear switch from 6506 to ethsw4948
- Moved management ethernet to Altix from ethsw02b to ethsw05a
- Convert all ISL's from leaf switches to 2-cable port channels
- Create link for submit via 10-gigabit

### **Node changes**

- Reimaged nodes and support machines with CentOS 4.5
  - Change smartd configuration so that it emails hpc-logs with drive SMART information. This will enable us to monitor drives that may be going bad and offline them prior to them going down unexpectedly.
- Upgraded RapidScale target and initiator software

- Upgrade Torque
- Upgrade Maui
- Drop the Topspin IB stack in favor of OFED-1.2
- Use OFED-provided MPI implementations (with tm interface for OpenMPI and other OFED annoyances fixed)
  - Upgrades for OpenMPI, MVAPICH, MVAPICH2
- Install OFED on the Altix
- Set system-wide default MPI implementations (OpenMPI/Intel)
- Configured ethernet only nodes to use ethernet binding to the RapidScale targets
- Retire iogw2, hpc, tp9400, osg
- Convert iogw4 to torque
- Convert hpcio9 to submit
- Install the following on submit:
  - NTP
  - DNS
  - LDAP
  - NFS
- Install the following on torque:
  - Torque
  - Maui
  - PBS log sweeper for website
    - perl-DBD-MySQL-2.9004-3.1.x86\_64.rpm
    - mysqlclient10-3.23.58-4.RHEL4.1.x86\_64.rpm
- Setup submit to be the submission node, which will then push all jobs over to Torque for the actual handling of these jobs.
- Rename all nodes to fit into new naming scheme of rack-side and slot number
- Upgrade to version 2.3.37 of LDAP
- Remove Pathscale compiler and software versions
- Remove older OpenMPI packages