OSG All Hands Meeting - Purdue Tier-2 Report
March 23rd, 2015

**Thomas Hacker**
for the Purdue Tier-2 Team
Erik Gough, Majid Arabgol, Manoj Jha, Norbert Neumeister

# Outline

- Site Resources
  - Computing
  - Storage
  - Networking
- Evolution of the Purdue site and campus developments that may affect Tier-2 site
- Completion of data center consolidation
- Monitoring

# Site Resources - Compute

- Purdue CMS compute resources fall into two groups
  - CMS owned Hadoop cluster
  - CMS owned nodes on Purdue's Community Clusters
- Hadoop cluster
  - Provides batch slots via Condor and HDFS storage for CMS
  - Moved to main campus (MATH) Datacenter in May 2014
- Community Cluster Program
  - Foundation of Purdue's research infrastructure
  - One cluster has been built per year with a common configuration
  - Shared by groups on campus who purchase nodes
  - Low overhead for research groups (no need to build/maintain your own cluster, just buy some nodes)
  - Node pricing on Community Clusters is very attractive due to volume

# 2015 Community Cluster

- Purdue Community Cluster program began in 2008
  - Installed six cluster systems over past six years
  - EDUCAUSE article July 2014 describes the program in detail
    > Hacker, T., Yang, B., McCartney G. (2014). Empowering Faculty: A Campus Cyberinfrastructure Strategy for Research Communities. EDUCAUSE Review Online, July 14, 2014.*

- The program has been very beneficial for Purdue and CMS
  - A portion of nodes were purchased with CMS funds
  - Substantial quantity discounts and institutional leveraging
  - The first (Steele) is now retired
  - Rossmann will be retired October, 2015
  - Hammer is Purdue's newest Cluster Community optimized for High Throughput Computing

* http://www.educause.edu/ero/article/empowering-faculty-campus-cyberinfrastructure-strategy-research-communities

# 2015 Community Cluster

- Purdue is fielding three cluster systems this year
  - Motivation: growing demand and need differentiation
  - High Throughput Computing, High Performance Computing, Need for very large memory nodes
  - High performance network technology (i.e. Infiniband) changes frequently enough to make it difficult to integrate large new clusters each year
- High Throughput Computing: *Hammer* and Big Memory nodes
  - Big memory nodes one system for life sciences
  - Primarily for jobs with no need for high-performance network (e.g. MPI jobs)
  - Will be expanded each year with new nodes
  - Old nodes will be retired from the system as they age out
  - Best suited for CMS needs; will trade-in old nodes
- HPC: *Rice*
  - Inifiniband for MPI
  - Same base nodes as Hammer system

# 2015 Community Cluster

- Each node in the Conte community cluster contains 2x8-core Intel Xeon E5 processors as well as an Intel PHI processor
  - Besides initial R&D, CMS has not used the Phi processors

- The 2015 cluster Hammer is focused on High Throughput Computing
  - We exchanged Conte nodes for Hammer nodes in 2015
  - Will gain an additional year of warranty
  - Purchasing additional Hammer nodes in 2015

- Current listing of CMS owned computational resources are shown in the next slide.
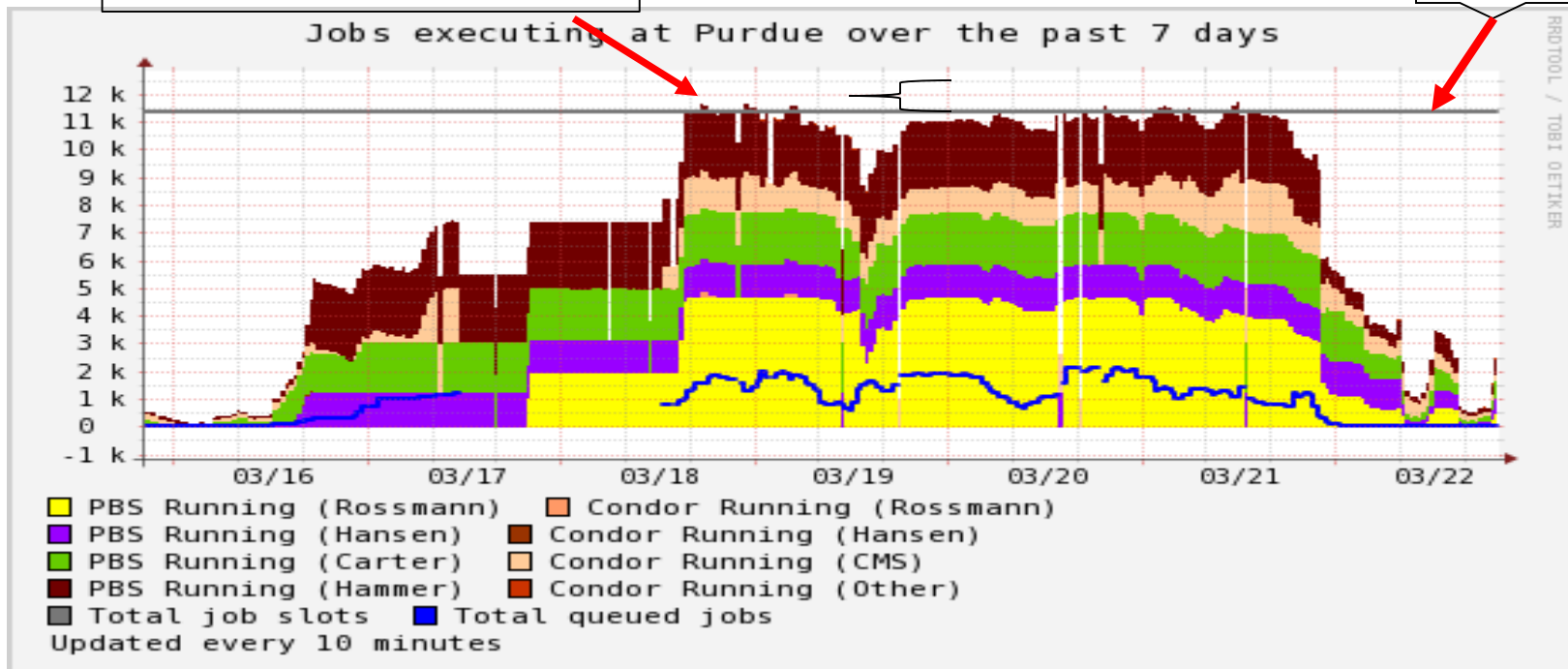
# Site Resources - Compute

| Cluster | Purchase Year | Max # nodes | Inter connect | #cores/node | Scheduler | Node Type | # CMS Nodes | Job Slots | #Years left |
|---|---|---|---|---|---|---|---|---|---|
| Purdue-Hadoop | Ongoing | 56 | 10 Gbps | 16 & 24 | Condor | Mixed | 56 | 1224 | Variable |
| Rossmann (to be retired Oct 2015) | 2010 | 438 | 10 Gbps | 24 | PBS/Condor | HP DL165 | 184 | 4632 | 0 |
| Hansen | 2011 | 201 | 10 Gbps | 48 | PBS/Condor | Dell R815 | 25 | 1200 | 1 |
| Carter | 2012 | 660 | 56 Gbps Infiniband | 16 | PBS | HP SL230 | 117 | 1872 | 2 |
| Hammer(-z) (temporary) | *2009* | *304* | *10 Gbps* | *8* | *PBS/Condor* | *HP Proliant* | *300* | *2432 (temp.)* | *To be replaced with new Hammer nodes this year* |
| Total | | | | | | | | 11,360 | |

# Site Resources - Compute

- CMS is able to utilize opportunistic slots available at Purdue. In the following diagram, opportunistic slots are represented by utilization above 11.3K line.
  - Total of 9,200 opportunistic slots available beyond 11.3K dedicated cores
  - Allows Purdue to reach running job counts above our available PBS queue capacity.  Max depends on the number of free slots in community cluster (PBS jobs have more priority and may evacuate Condor jobs if resource is needed).

# Site Resources - Storage

- File systems
  - Hadoop Distributed Filesystem (HDFS)
    - Mixture of Compute/Storage (running Condor) and pure storage (Hadoop data nodes)
  - SCRATCH area
    - Each community cluster has highly scalable parallel file system (currently Lustre). User can use this area for reading and writing file from/to this area. 500 GB of space is allocated to user in their scratch area.
  - HOME area
    - Isilon system at Purdue provides HOME area on all the worker nodes at the site. Up to 100 GB of space is allocated to user's in their home area.
  - Group space
    - Migrating from old system to new Purdue Data Depot
    - Higher performance than old system
  - Common storage resources (i.e. HOME) are provided by central computing at Purdue (ITaP) at no charge to the project

# Site Resources - Storage

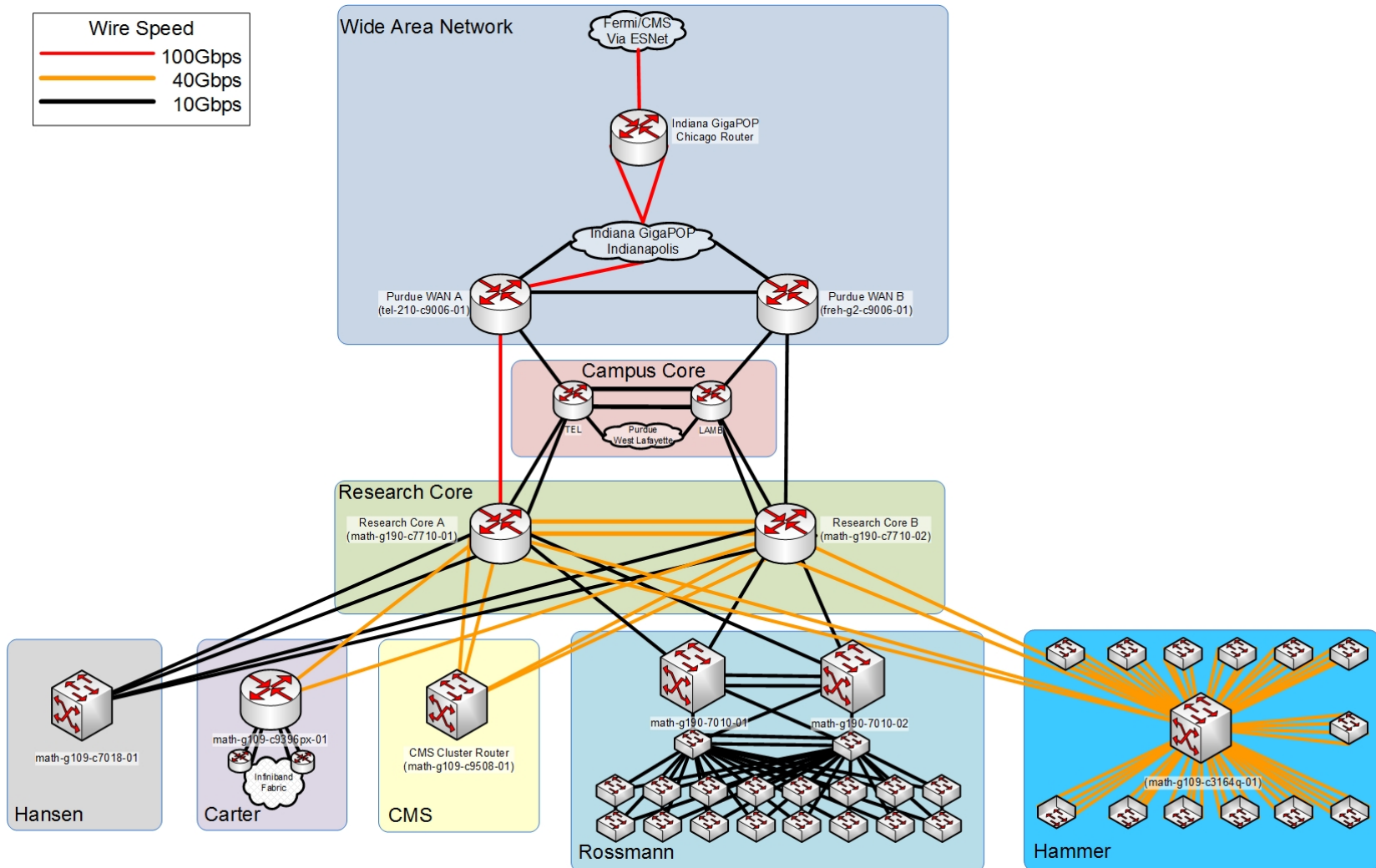| Node Type | # of Nodes | Raw Disk (TB) | Total Raw Disk (TB) |
|---|---|---|---|
| HP DL165 | 41 | 8 | 328 |
| Kingstar (1st gen.) | 13 | 48 | 624 |
| Kingstar (2nd gen.) | 10 | 108 | 1080 |
| Kingstar (3rd gen.) | 2 | 64 | 128 |
| Advanced HPC | 7 | 48 | 336 |
| Advanced HPC | 4 | 144 | 576 |
| Advanced HPC | 5 | 144 | 720 |
| Kingstar (4rd gen.) | 7 | 144 | 1008 |
| Dell XD 720 | 8 | 48 | 384 |
| | | **Totals** | **5184 TB Raw** <br> **4696 TiB Raw** <br> **~2.2 PB Usable @ 2 reps** |

# Hadoop Updates in 2014

- Patched OSG HDFS release to prevent timeouts during block reports on large datanodes (144TB, 2m+ blocks)
  - Released in OSG HDFS version 2.0.0+545-1.cdh4.1.1.p0.20
- Disabled Transparent hugepage compaction
  - During xrootd scale testing, we observed high system cpu utilization on nodes running both HDFS and gridftp
  - `echo never > /sys/kernel/mm/redhat_transparent_hugepage/defrag`
- Disabled zone reclaim mode
  - During 20 Gbps testing, more system cpu spikes were discovered on gridftp/HDFS servers
  - `echo 0 > /proc/sys/vm/zone_reclaim_mode`

# Site Resources - Networking

- Changes from prior year marked in <span style="color:red">RED</span>
- Wide Area Networking
  - WAN link to Chicago: <span style="color:red">100 Gb/sec Connectivity to FNAL</span>: CMS resources at site connect to FermiLab through Indiana GigaPop Chicago router, and then via <span style="color:red">Fermi/CMS via Esnet</span>
  - <span style="color:red">Bandwidth between Indiana GigaPop and Fermi/ESnet was changed from a 10G dedicated to a 100 Gbps shared link</span>
- Data Center Core and Local Area Networking
  - Community Clusters to Data Center Core: 40 Gbps
  - <span style="color:red">CMS Cluster to Data Center Core: 160 Gbps</span>
  - CMS Cluster nodes: 10G Ethernet
  - Community clusters: 10G Ethernet or 40G/56G Infiniband
- Project funds high speed WAN connectivity

# Current Networking Map

# Status of 2014 Technology Goals
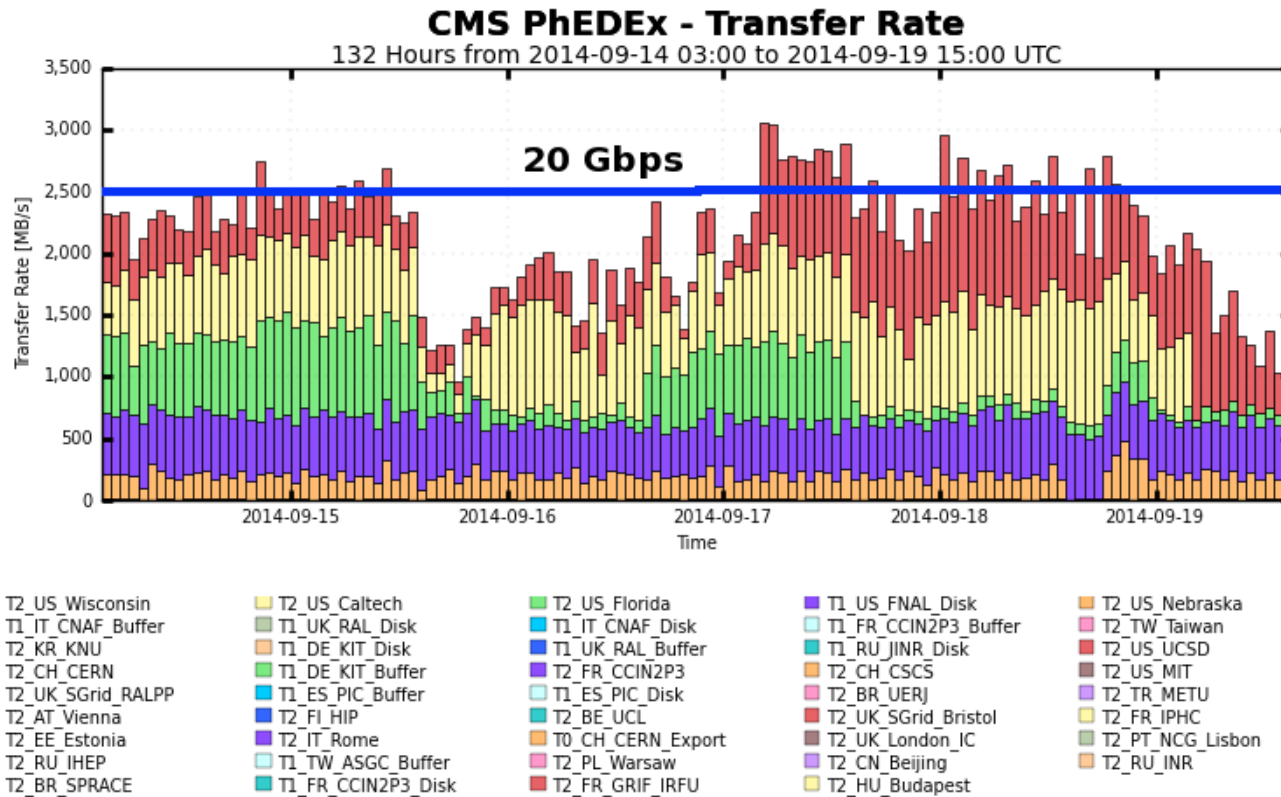
- Completed
  - OSG 3.2 Upgrade
  - HTCondor-CE (one CE operational)
  - HTCondor 8.2.7
  - Config Management (puppet)
  - 20 Gbps test
  - Xrootd 2k/4k tests
- Partially Completed
  - Services IPv6
  - WNs IPv6

# IPv6 Status

- Goal is for all CMS servers and Community Cluster WNs to be dual stacked IPv4/IPv6
- Worker Nodes with IPv6 connectivity
  - CMS Cluster: All nodes enabled
  - Community Clusters: Rossmann and Hammer
- Services with IPv6 enabled
  - Xrootd, gridftp servers
- TODO
  - Enable IPv6 on SRM
  - Enable IPv6 on Hansen and Carter

# 20 Gbps Test

- Completed 20 Gbps test 9/17/14



CMS PhEDEx - Transfer Rate
132 Hours from 2014-09-14 03:00 to 2014-09-19 15:00 UTC

# Evolution of Purdue Tier-2 Center

- Datacenter Consolidation
- Networking Updates
  - LHCONE Peering enabled
  - IPv6
- Moving to the Hammer HTC cluster
- Local Initiatives
  - Monitoring
  - Gatekeeper experiences

# Purdue CMS Hardware

# Datacenter Consolidation

- CMS Datacenter Consolidation
  - In May 2014, CMS owned cluster and HDFS storage moved to main Research Computing datacenter in Math building
- Benefits for CMS
  - New networking equipment, all nodes upgraded to 10Gbps
  - Higher core networking bandwidth
  - No cost replacement of out of warranty storage/compute nodes
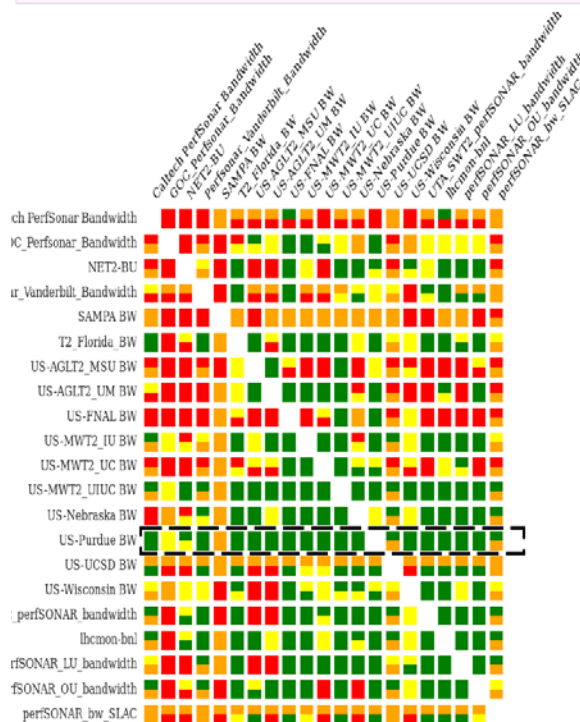
# LHCONE Peering

- Network traffic at our site integrated to LHCONE.

- Observing improved connectivity with other Tier-1 and Tier-2s.

- Around 10K of single core job slots, Purdue Tier-2 site is in better place to efficiently utilize data using dynamic data caching and AAA mechanisms.



**OSG Grid Operations Center Test Mesh Config Dashboard**

OSG Grid Operations Center Test Mesh Config - GOC Bandwidth Tests

Throughput >= 900Mbps | Throughput < 900Mbps | Throughput <= 500Mbps | Unable to retrieve data | Check has not yet run

# Multicore Scheduling

- Multicore glideins are used to run both production and analysis jobs for PBS clusters at the site

  – Each glidein requests for 8 single cores per job.

- Observed batch scheduling efficiency increased after switching to multicore glideins.

- Sometimes, we observed number of unclaimed glideins varies between ~ 5 to 30% of total glideins running at the site.

# Gatekeeper reduction

- Each community cluster at Purdue needs its own computing element for gatekeeper to OSG
- Currently, we have 5 clusters and 6 gatekeepers
- Disadvantages
  - Lots of CEs to monitor and manage
  - Some difficulties in load balancing across clusters
- Reduction of gatekeepers over time
  - Each year as we acquire Hammer HTC cluster nodes and retire older nodes
  - Retire old gatekeeper with old cluster nodes
- Eventually, we plan to operate two to three gatekeepers to a single large CMS queue

# Local Monitoring

- Sensu deployed across all CMS servers and community cluster WNs
  - Monitor proxy cert expiration, keepalive, external WAN connectivity, SSH, PBS mom on WN and scheduler, kernel version, disk space
  - Alerts are sent when a test result is critical
  - GUI for checking overall status
- Developed home grown scripts for monitoring Hadoop cluster
  - Disk space, stale processes, gridftp xfers, network bandwidth
  - Problematic nodes change their color from green to red
  - http://web.rcac.purdue.edu/cms/Workers/Services/Hadoop/table_hosts.html

# User Support

- Providing disk space to `store/user' and login access to user interface for users from Carnegie-Melon, Ohio State, Purdue and SUNY-Buffalo Universities.
- Users get prompt response in case of any problem related to utilization of the resources at our site.
- Analysis datasets from different physics groups are managed through Dynamic Data Management.  However, allotted some dedicated storage disk space for different physics groups at site.

| Group | Subscribed | Resident |
|-------|-----------|----------|
| AnalysisOps | 717.80 TB | 713.72 TB |
| exotica | 32.18 TB | 31.18 TB |
| heavy-ions | 7.03 TB | 7.03 TB |
| upgrade | 16.52 TB | 16.52 TB |
| b-physics | 329.81GB | 329.81 GB |
| higgs | 4.20 GB | 4.20 GB |
| Local | 92.48 TB | 92.48 TB |

# Summary and Preparation for Run 2

- Work in 2014 in preparation for Run 2
- Improved network connectivity to/from site
  - 100 Gb/sec path to FNAL
  - LHCONE connectivity
- 10k+ single core job slots
- Deployed multicore glideins for both production and analysis jobs on PBS clusters
- Purdue Tier-2 site is well positioned to efficiently utilize data using dynamic data caching and AAA mechanisms.

# Acknowledgements

- Thanks to the Purdue CMS Tier-2 team for their efforts and contributions to this talk
  - Erik Gough
  - Majid Arabgol
  - Manoj Jha
- Thanks to ITaP
  - John Wright (past)
  - Nikki Huang (current)

# Questions?