

# MIT Tier2 site report



Max Goncharov  
March 8, 2010

B.Alver M.Betancourt S.Jaditz W.Li C.Paus D.Ross G.Stephans M.Tiernan  
L.Winslow B.Wyslouch

USCMS Tier2 workshop

- Hardware Infrastructure
- Move to the New Location
- Software Development
- Monitoring
- Performance
- Future Plans

2x4-core Intel 2.3 GHz 3.75 TB	x20
2x4-core AMD 1.6-2.0 GHz 2.8 TB	x70
2x4-core AMD 2.0 GHz 1.86 TB	x55
2x2-core AMD 1.8 GHz 1.6 TB	x20
Tier2 funded	

2x2-core AMD 2.0 GHz 1.1 TB	x30
2x1-core AMD 1.6 GHz	x200
Other resources (CDF/CMS-HI)	

## Location: Bates

- Hosted by MIT, 22 mi from Cambridge
- Nodes – mostly Dell/Thinkmate 2U multiple core with 6 or 8 disks for storage
- Dell PowerEdge 2950 and R710
- Storage – dCache on individual nodes with software RAID 5

## CPU power and batch slots:

Tier2 resources	Batch slots	1270
	CPU power	9252 HS06
CDF/HI	Batch slots	530
	CPU power	1728 HS06
Total	Batch slots	1800
	CPU power	10980 HS06

## Storage figures:

Tier2 resources	395 TB
Other (CDF/HI)	28 TB
Total	423 TB



# Move to New Location - Bates

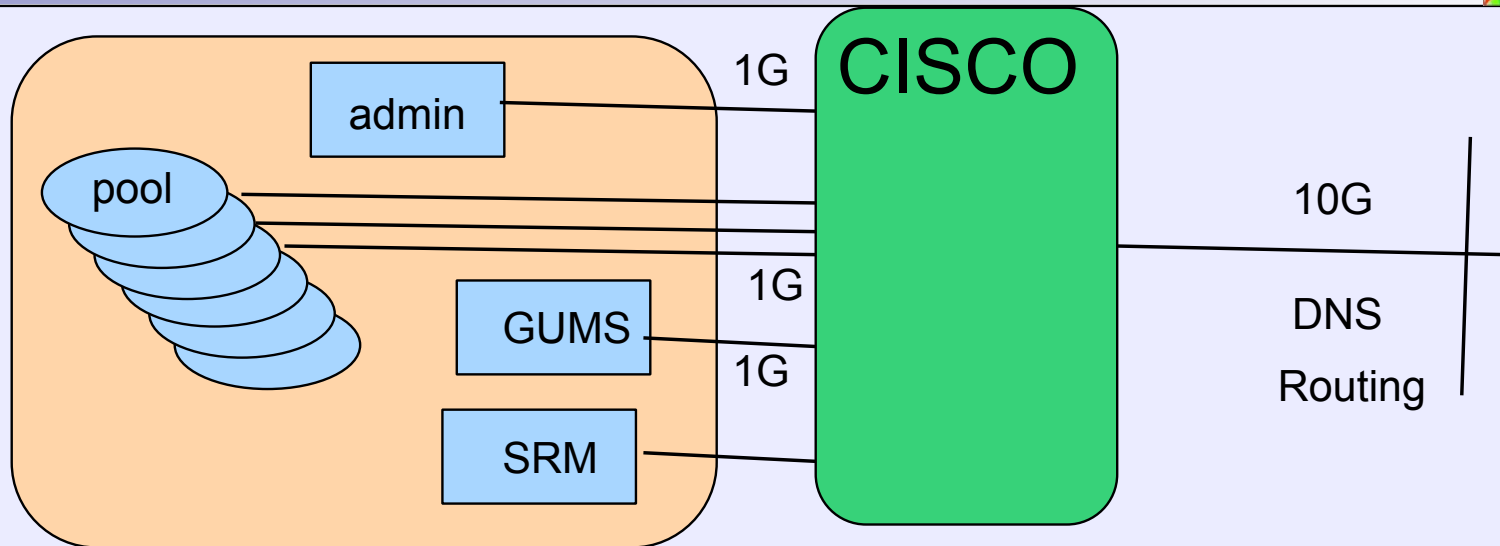


## Last Year:

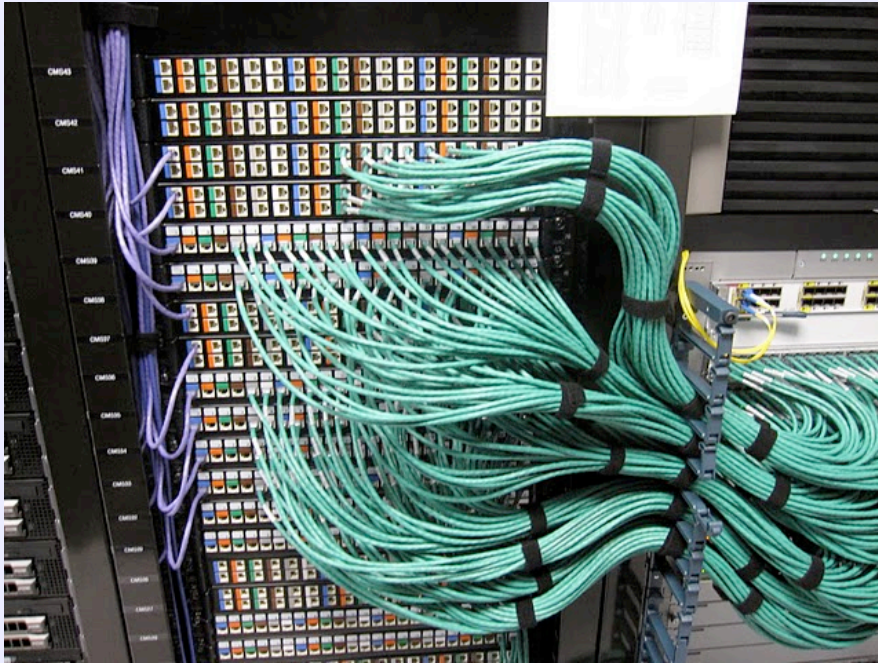
- location hosted by MIT IS&T, on campus
- overseen 5x16 by operators
- total size – 15 racks

## Nov 2009-February 2010:

- first we finalized the public network setup
- than moved to Bates, location hosted by MIT, 22 mi from Cambridge
- network managed by MIT IS&T
- racks, power, infrastructure managed by Bates
- overseen 7x24
- UPS backup for all servers (that are not worker nodes) (4 racks)
- 30 water cooled racks, rack – 40 U and 10 kW
- declared downtime on Monday, were running jobs on Friday



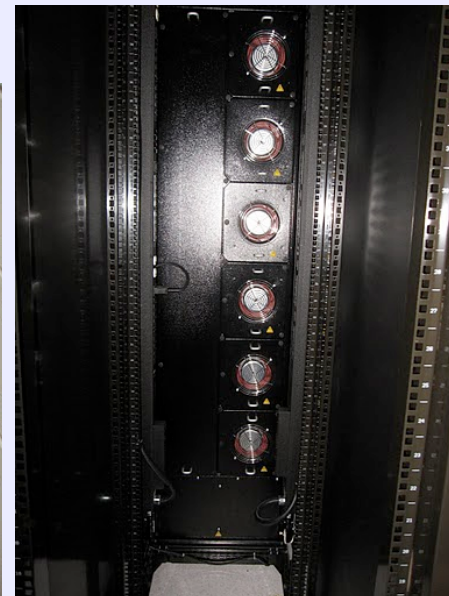
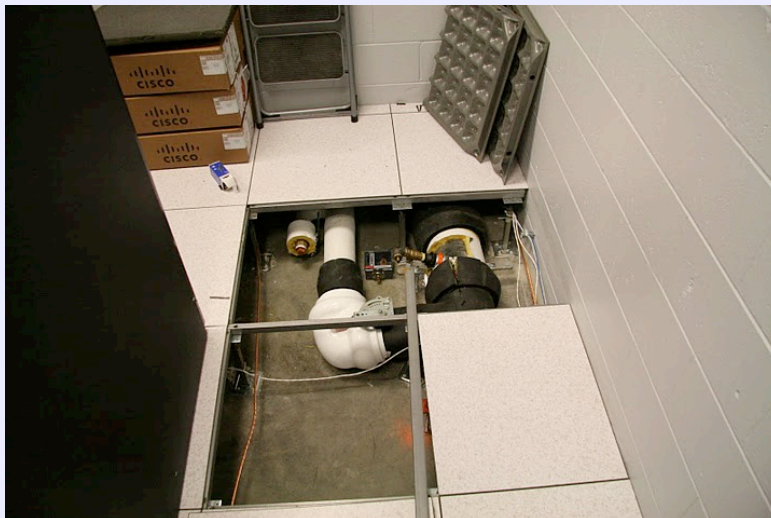
- Leased for 5 years CISCO Nexus 7016
- Machines can talk at 1Gb through copper cable links
- Before moving to Bates switched to public network
  - all IPs changed, kept same DNS names
  - reconfigured ALL services
  - developed firewall rules
- dCache – each pool now serves as gridFTP
- Declared downtime on Monday, were running jobs by Thursday



Network is huge improvement over the last arrangement

Water Cooled Racks

- 40 U, 10 kW of power



## This year, 2010, at the same time as Bates move

- Bought and installed 20 R710 Dell nodes
  - 2x4core CPUs (24GB mem) , Nehalem → 16 Condor slots max (running 12)
  - 6 1TB disks → 4 TB available for dCache
- Bought and installed 2 R710 Dells
  - one is Compute Element
  - second one is dCache admin
- Bought and installed R710 Dell with 15 K disks
  - our new NFS server (osg, app, opt)
  - finished switching services last month

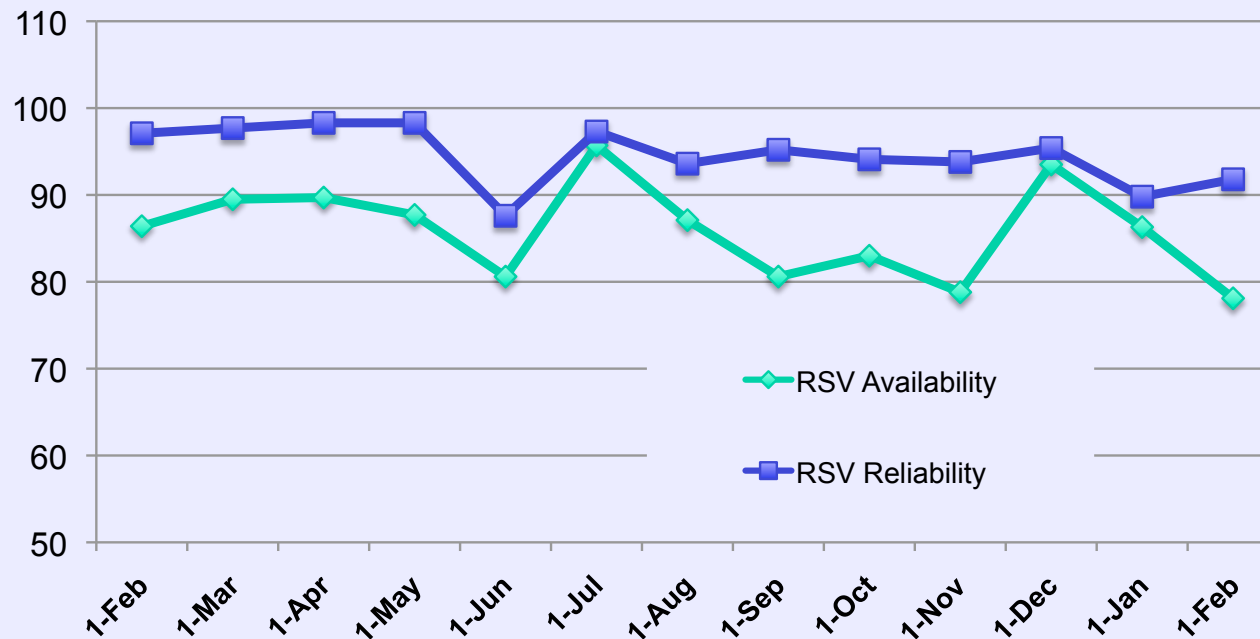


- Tested using hardware RAID instead of software one
  - system disk RAIDed, no single disk failure downtime
  - better utilize space for dCache
    - 3.7 TB available with software RAID
    - 4.2 TB available with hardware RAID on the same node
  - improvement in read/write access
- We plan to start using hardware RAID on newly purchased systems

- OSG: switched from version 0.8 to 1.2
- Condor: switched from version 6.8 to 7.2
  - adjusted scheduling to round robin to alleviate pre-emption
  - wrote custom monitoring to prevent “black holes” nodes
- dCache: switched from version 1.8 to 1.9.3
  - made all pools gridFTP doors
  - installed host certificates for each dCache pool
- PhEDEx: switched from ... to 3\_2\_9
- Switched from SL4 to SL5 all worker nodes, new servers are installed with SL5
- Deployed central NFS server
  - hosts all common areas
  - simplified "CA/CRL" certificate synchronization

## Nagios:

- Re-worked monitoring of individual machines
  - each host has 10-15 checks
- Added AIM/e-mail into Nagios, will be promptly notified if
  - if any node goes down, running without firewall
  - if disk fails out of raid
- Added cron scripts to monitor and take action if
  - “black hole” node appears
  - dCache stops functioning (dccb/srm not working)
  - SAM/RSV pages show problems
- Automated response to facility problems
  - shut cluster down if temperature goes up to dangerous levels



Two long downtime due to the move to Bates  
 Bates is new facility:

- long downtime due to broken UPS support
- long downtime due to the Chiller quitting in cold weather

## Hardware

- Add 200 TB to dCache space
  - buy Dell servers with 2/1 TB disks and 2x4 cores
- Investigate other vendors and other storage options

## Monitoring

- There is plenty of information available to monitor conditions at Bates
  - all racks provide fan speed, water flow, temperature, ...
  - UPS condition are available
- Improve our automated response to problems

- MIT Tier-2 is now in its final location
- Plenty of room for expansion
- We have significantly improved monitoring and operating procedures
- We maintained reasonable performance while undergoing many changes
- We are working to have higher availability and reliability in the months to come



# Acknowledgements



We received tremendous help from many people

Tier2: Ken Bloom, Bockjoo Kim, Stefano Belforte, Arvind Gopu and others

Fermilab Team : Burt Holzman, Paul Rossman, Brian Bockelman, Catalin Dimitresku, Denis Perelmutov, John Weigand, Anthony Tiradani and others

Thanks to UCSD team for guidance during the network switch