



# Florida Tier2 Site Report

Presented by Yu Fu

for the University of Florida Tier2 Team

(Paul Avery, Bourilkov Dimitri, Yu Fu, Bockjoo Kim, Yujun Wu)

USCMS Tier2 Workshop

Livingston, LA

March 3, 2009

# Outline

- Site Status
  - Hardware
  - Software
  - Network
- FY2009 Plan
- Experience with Metrics
- Other Issues



# Hardware Status

- Computing Resources:
  - UFlorida-PG (merged from previous UFlorida-PG and UFlorida-IHEPA):
    - 126 worknodes, 504 cores (slots)
    - 84 \* dual dual-core Opteron 275 2.2GHz + 42 \* dual dual-core Opteron 280 2.4 GHz, 6GB RAM, 2x250(500) GB HDD
    - 794 kSI2k, RAM: 1.5 GB/slot
    - Older than 3 years, out of warranty, considering the possibility to upgrade or replace sometime in future.

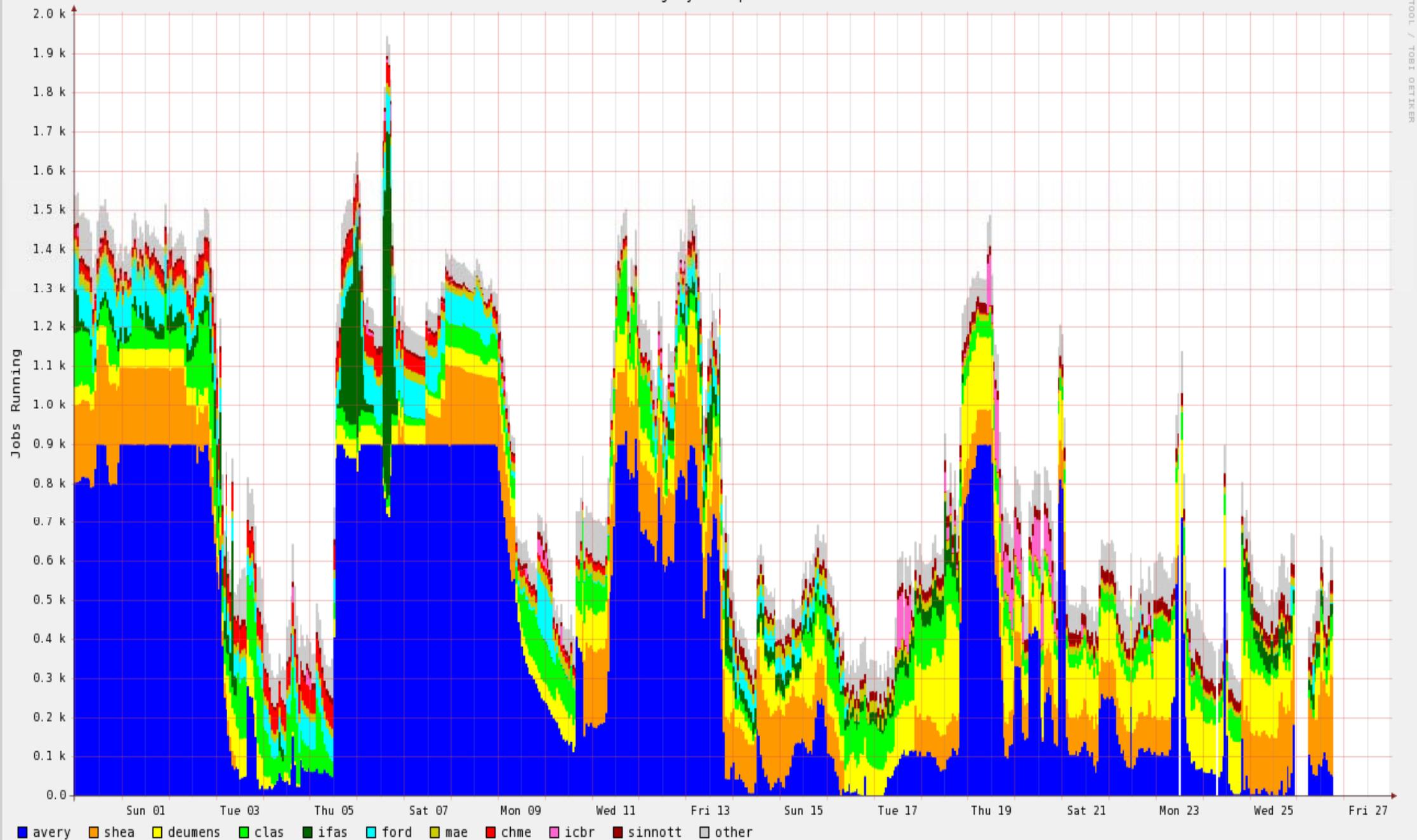
# Hardware Status

## – UFlorida-HPC:

- 530 worknodes, 2568 cores (slots)
- 112 \* dual quad-core Xeon E5462 2.8GHz + 418 \* dual dual-core Opteron 275 2.2GHz
- 5240 kSI2k, RAM: 4GB/slot, 2GB/slot, 1GB/slot
- Managed by UF HPC Center, Florida Tier2 invested partially in three phases.
- Tier2's official share/quota is 900 slots (35% of total slots), and Tier2 can use more slots on opportunistic basis. The actual average Tier2 usage is ~50%.
- Tier2's dedicated SpecInt: 1836 kSI2k

# HPC cluster usage of last month

Jobs Running by Group (last month)



# Hardware Status

- **Interactive analysis cluster for CMS**
  - 5 nodes + 1 NIS server + 1 NFS server
  - 1 \* dual quad-core Xeon E5430 + 4 \* dual dual-core Opteron 275 2.2GHz, 2GB RAM/core, 18 TB total disk.
- **Total Florida CMS Tier2 dedicated computing power (Grid only, not including the interactive analysis cluster):**  
2.63M SpecInt2000, 1404 batch slots (cores).
- **Have fulfilled the 2009 milestone of 1.5M SpecInt2000.**
- **Still considering to get more computing power in FY09.**

# Hardware Status

- Storage Resources:
  - Data RAID: gatoraid1, gatoraid2, storing CMS software, \$DATA, \$APP, etc. 3ware controller with SATA drives, mounted as NFS. Very reliable.
  - Resilient dCache: 2 x 250 (500) GB SATA drives on each worknode. Acceptable reliability, a few failures.
  - Non-resilient RAID dCache: FibreChannel RAID (pool03, pool04, pool05, pool06) + 3ware-based SATA RAID (pool01, pool02), with 10GbE or bonded 1GbE network. Very reliable.
  - HPC Lustre storage, accessible both directly and via dCache.
  - 20 GridFTP doors: 20x1Gbps

# Hardware Status

– Total dCache Storage:

Resource	Raw	Usable	Usable (after 1/2 factor for resilient)
pool01	9.6TB	6.9TB	6.9TB
pool02	9.6TB	6.9TB	6.9TB
pool03	18.0TB	13.9TB	13.9TB
pool04	18.0TB	13.9TB	13.9TB
pool05	54.0TB	47.2TB	47.2TB
pool06	72.0TB	55.2TB	55.2TB
worknodes	93.0TB	71.1TB	35.6TB
HPC Lustre	35TB	30TB	30TB (space actually used by T2)
Total	309TB	245TB	210TB

(Hard drives in the UFlorida-HPC worknodes are not counted since they are not deployed in dCache system.)

# Hardware Status

- Total raw storage is 309TB, real usable Tier2 dCache space is 210TB, still some gap to meet the 400TB target of '09.
- Planning to deploy 200TB new RAID in 2009.

# Software Status

- Most systems running 64-bit SLC4, some have migrated to SLC5
- OSG 1.0.0
- Condor (UFlorida-PG and UFlorida-IHEPA) and PBS (UFlorida-HPC) batching system.
- dCache 1.9.0
- Phedex 3.1.2
- Squid 3.0
- GUMS 1.2.16
- .....
- All resources managed with a 64-bit customized ROCKS 4, all rpm's and kernels are upgraded to current SLC4 versions.
- Preparing ROCKS 5.

# Network Status

- Cisco 6509 switch
  - All 9 slots populated
  - 2 blades of 4 x 10 GigE ports each
  - 6 blades of 48 x 1 GigE ports each
- 20 Gbps uplink to campus research network
- 20 Gbps to UF HPC
- 10 Gbps to UltraLight via FLR and NLR
- Florida Tier2's own domain and DNS
- All nodes including worknodes are on public IP, directly connected to out world without NAT.
- UFlorida-HPC and HPC Lustre with InfiniBand.

# FY09 Hardware Deployment Plan

- Computing resources
  - Have already met the 1.5M SI2k milestone, still considering to add more computing power.
  - Investigating two options:
    - Purchase a new cluster: power, cooling and network impact.
    - Upgrade present UFlorida-PG cluster to dual quad-core: no additional power and cooling impact, may be more cost-effective but re-used old parts may be unreliable.
  - No official SpecInt2000 numbers available for new processors?

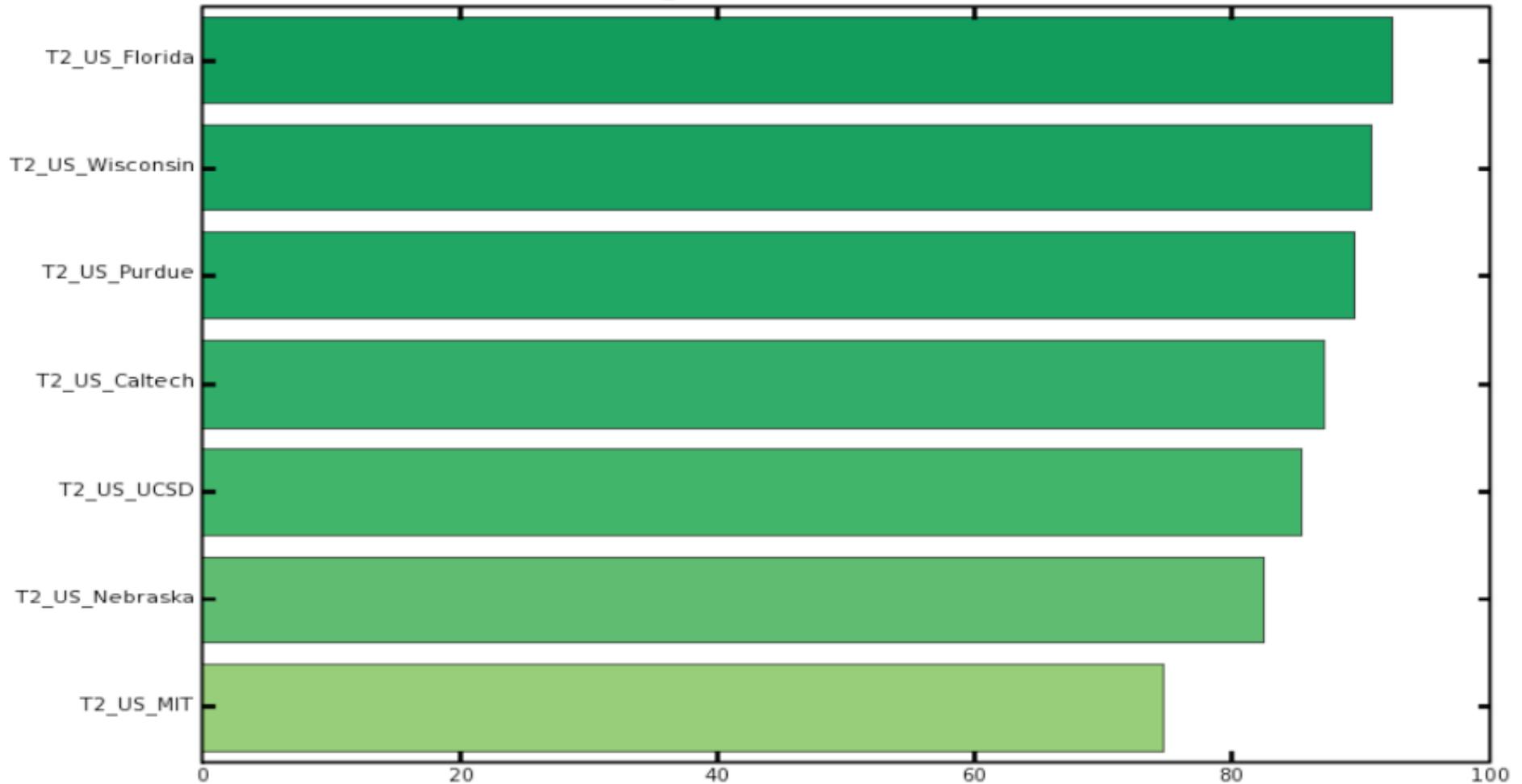
# FY09 Hardware Deployment Plan

- Storage resources
  - Will purchase 200TB new RAID.
  - To avoid network bottleneck, don't want to put too much disk under a single I/O node.
  - Considering 4U 24-drive servers with internal hardware RAID controllers.
  - Waiting until 2TB drives become reasonably available at enterprise-level stability: 1TB drives require more servers and they will be made obsolete soon by the 2TB ones.

# Experience with Metrics

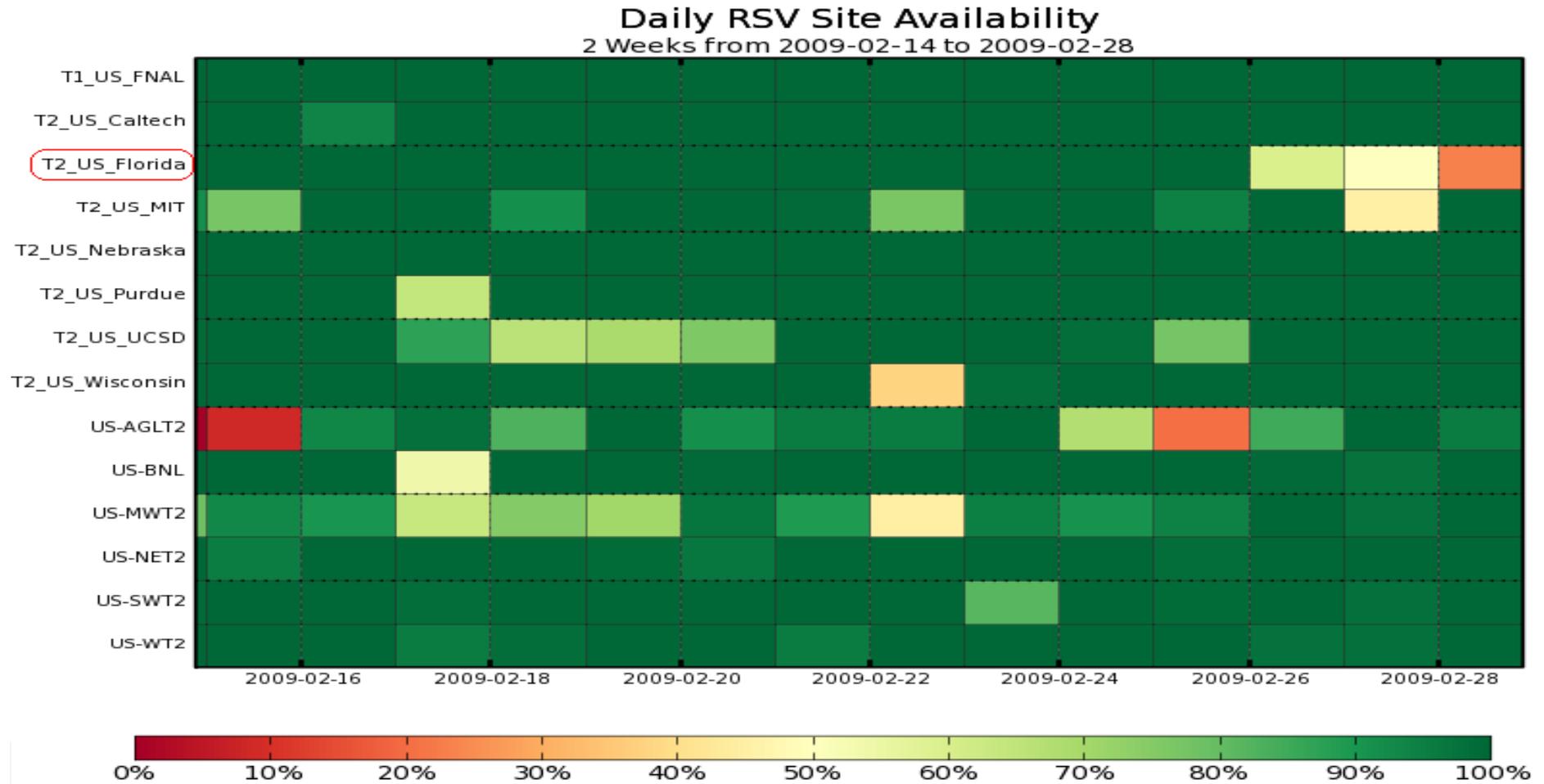
- SAM: excellent

**Site Availability, 2008-09-01 - 2009-03-01**



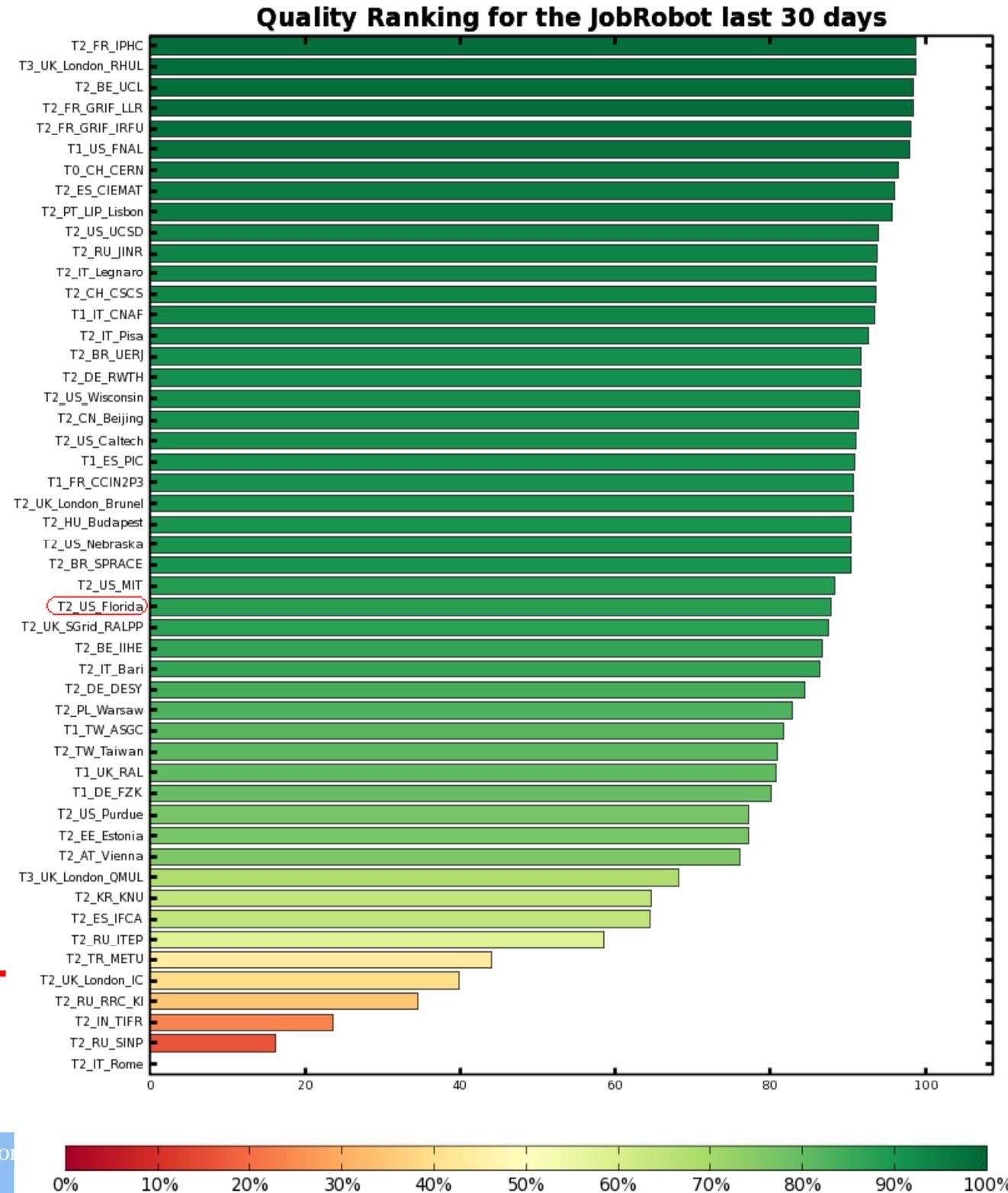
# Experience with Metrics

- RVS: good



# Experience with Metrics

- JobRobot: OK
  - We had three (now two) different clusters.
  - Limited available slots on small clusters – merging may help.
  - Proxy expiration problems.
  - Ambiguous errors due to glite job submission.



# Experience with Metrics

- We have self monitor to monitor SAM, RVS, JobRobot metrics monitoring systems.
- Our monitor systems notify us the problems instantly by emails – this always allows us to fix problems as quickly as possible.
- The alarm emails also help us think what tools we need to develop and improve to better monitor, diagnose and fix the problems.
- The tools developed with alarm emails have proved to be very useful.

## Florida : Global Monitors

SAM CMS	<a href="#">SAM Test Page</a>	<a href="#">Monitoring SAM</a>	<a href="#">Site Status</a>		
JobRobot	<a href="#">JobRobot Page</a>	<a href="#">S.E.=99%</a>	<a href="#">Site Status</a>		
CMSSW	<a href="#">ALL</a>	<a href="#">T2+T3</a>	<a href="#">OSG T2+T3</a>	<a href="#">CsDoGrid</a>	<a href="#">OSG Twiki</a> <a href="#">GridCat</a>
Dataset	<a href="#">Requests: View1 View2</a>	<a href="#">Arrivals</a>	<a href="#">Datasets at Florida</a>	<a href="#">Data MC</a>	
CMS Jobs	<a href="#">Job Summary</a>	<a href="#">PG</a>	<a href="#">IHEPA</a>	<a href="#">HPC</a>	
SAM OSG	<a href="#">SAM OSG-CE Test</a> <a href="#">SAM OSG-SRMv2 Test</a>	<a href="#">CE</a> <a href="#">SRMs?</a>	<a href="#">Check Log</a>	<a href="#">Gratia Report Page</a>	<a href="#">CEMon/BDII Page</a>
dCache INFO	<a href="#">dCache</a>	<a href="#">cell status</a>	<a href="#">pool status</a>	<a href="#">srmwatch</a>	
PhEDEX Prod.	<a href="#">Download Quality</a> <a href="#">Upload Quality</a>	<a href="#">PhEDEX Rate</a>	<a href="#">PhEDEX Status</a>	<a href="#">Request</a>	<a href="#">C.W.O.T2</a>
PhEDEX Debug	<a href="#">Download Quality</a> <a href="#">Upload Quality</a>	<a href="#">PhEDEX Rate</a>	<a href="#">PhEDEX Status</a>		<a href="#">C.W.O.T2</a>

## Florida Clusters : Job Monitor

<a href="#">UFlorida-PG</a>		<a href="#">CMSSW Installed</a>	<a href="#">Ganglia</a>			
Account Name	(Real Name)	Idle	Running	Held	Other	Total
cms30010	(xxxxxxx...)	0		0	1	1
cms30477	(xxxxxxx...)	0		0	1	1
cms30637	(cardaci/CN=6564...)	0		22	0	22
cms30672	(xxxxxxx...)	0		0	17	17

# Experience with Metrics

- We also have various other monitor systems: Ganglia, Nagios and systems to monitor all aspects of hardware, software as well as services like PhEDEx, dCache, Squid, status of dataset transfer, ..... etc.
- Operations support helps only if a useful solution, suggestion or hint is provided.
- We often find we have to understand the details of what operations support is doing, this can take quite some time.

# Other Issues

- Merge UFlorida-PG and UFlorida-IHEPA. Working on a prototype Condor-PBS combined gatekeeper based on random selection.
- Overloaded gatekeepers – to upgrade the hardware.
- What is the best way to deal with out-of-warranty old machines?
  - Becoming unstable and unreliable.
  - Parts can be expensive and hard to find.
  - Significantly lower performance or capacity than new ones - less energy-efficient.

# Other Issues

– Question: continue to maintain, de-commission, or upgrade?

– An example:

We considered to upgrade UFlorida-PG's gatekeeper's memory to 16GB, it turned out that 8\*2GB DDR memory (registered ECC) would cost more than \$1000. Finally we decided to upgrade the motherboard, processors and memory, the total cost is less than \$1000, but we got:

- New motherboard with better chipset
- New more efficient heatpipe-type heatsinks
- Dual dual-core processors -> Dual quad-core faster processors
- 4GB DDR RAM -> 16GB DDR2 RAM (faster than 16GB DDR)

– Bottom line: it can be more cost-effective to get a new machine or upgrade to a new machine than to fix an old one.