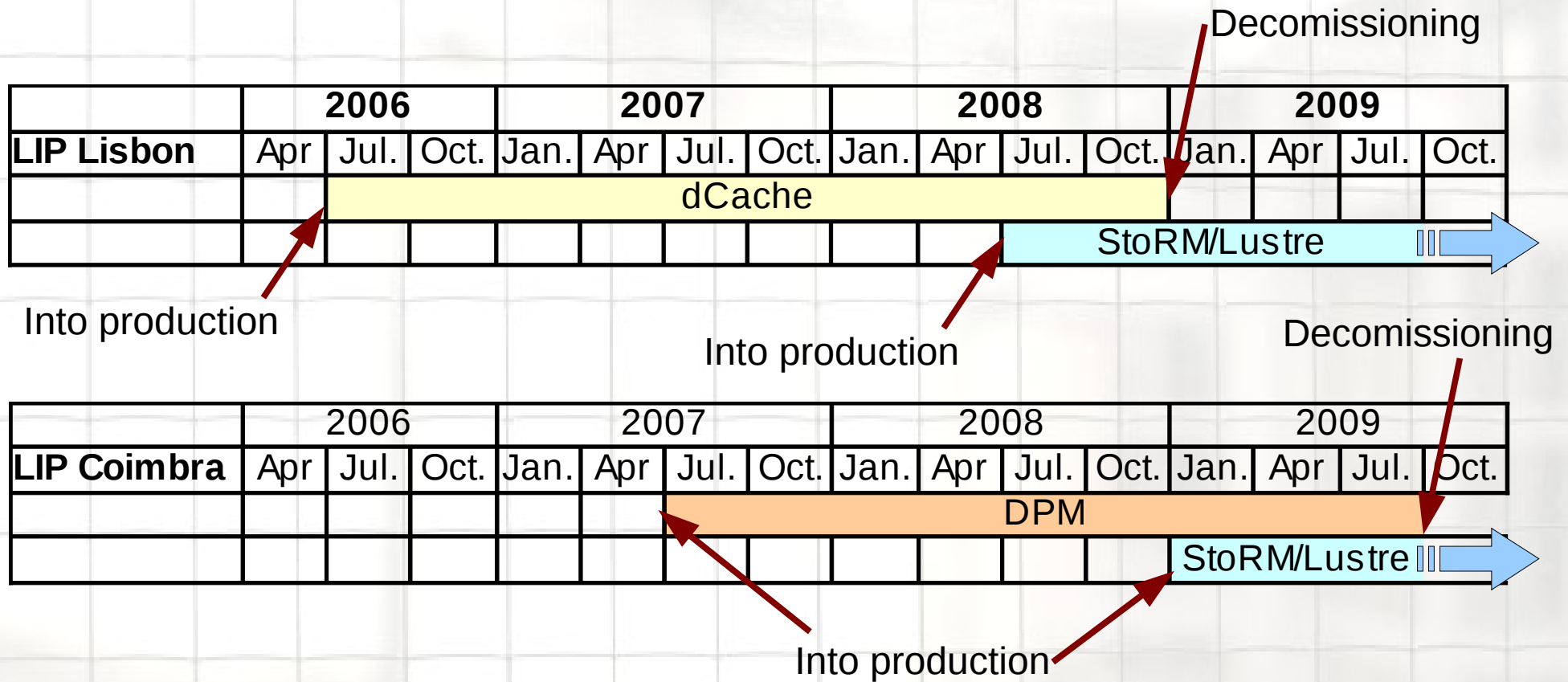


Operation of the Portuguese Tier2 (Storage Element component)

M. David, M. Oliveira

G. Borges, J. Gomes, J.P. Martins



- **dCache@Lisbon:**
 - **SE choice:**
 - **Could be used by Grid and local users.**
 - **Issues:**
 - **Needs too many human and computing resources to be adequate for a provider like LIP.**

- **DPM@Coimbra:**
 - **No need to support local users in the same SE system.**
 - **Rather simple to setup and administer.**
 - **Robust and reliable.**

The StoRM SRM with Lustre FS

- Middle of 2008: New storage and computing resources were purchased.
- A decision to change the SE technology was taken: Testing and deployment of Lustre FS from SUN and StoRM SRM both at the Lisbon and Coimbra sites.
- Lisbon, Lustre (version 1.6.6) servers:
 - **1 MGS + several MDT's:**
 - 1 HP Blade, FS's on iSCSI DS.
 - **5 OSS' s Dell PE1950:**
 - 20 OST's 130TB (Atlas + CMS)
 - **Other VO's and local users (~20TB):**
 - 1 OSS with 1 OST HP
 - 4 OSS's, older machines.
- StoRM (version 1.3.20):
 - **StoRM frontend and backend:**
 - 1 Dell PE1950
 - **2 GridFTP servers.**
 - 1 Dell PE1950 and 1 SUN Fire X2100
- Coimbra Lustre (version 1.6.6) servers
 - **1 MGS + several MDT's:**
 - 1 SUN Fire X2100
 - **2 OSS's Dell PE1950:**
 - 12 OST's 75TB (Atlas)
 - **Local users (12TB):**
 - 2 OSS's, older machines
- StoRM (version 1.3.20):
 - **StoRM frontend and backend in separate machines:**
 - 2 SUN Fire X2100
 - **2 GridFTP servers.**
 - 2 SUN Fire X2100

- StoRM decouples the SRM implementation from the underlying FS.
- The FS can be used by Grid and local users.
- Smooth learning curve for both the Lustre FS and StoRM.
- After ~9 months in operation it proved stable and robust.
- Clear “logs” for both Lustre and StoRM.
- Easy implementation of a backup policy.
 - **Backup and recovery of the namespace filesystems (MDT's) was successfully tested.**

- Lustre 1.8.0: released May 6:
 - **Needed features:**
 - Pools/Groups of disks.

- StoRM 1.4.0: released May 15:
 - **Needed features:**
 - Space Tokens limitation (at StoRM BE configuration level).
 - Highly configurable space area policies.
 - Much improved Administration Manual (*I have seen it!!*).
 - Real dynamic information provider (is a cron job in 1.3.20).

- We will deploy and test them soon.

- **Lustre:**
 - **Lagging behind in version of the supported kernels.**
 - **Communications between clients and servers → timeout problems → hanging clients**
 - **Solved on the network side**

- **StoRM:**
 - **Install/config too tied to the Italian Grid specificities (work in progress to integrate in the main gLite stack).**
 - **Still missing the 64 bit version (though work is in progress).**

- **Lisbon:**
 - **Nagios: problem detection and notification**
 - **Cacti: detailed network graphs**
 - **Ganglia: detailed graphs and aggregated graphs**
- **Coimbra**
 - **Nagios + pnp: problem detection and notification and detailed graphs.**

Nagios@Lisbon

se003	/	OK	05-15-2009 08:54:45	51d 15h 12m 16s	1/3	/: 40%used(8572MB/21542MB) (<90%) : OK
	eth0	OK	05-15-2009 08:56:01	51d 20h 44m 8s	1/3	OK: eth0: Link is up at 1000Mb/s, Full duplex
	eth1	OK	05-15-2009 08:56:03	51d 20h 42m 35s	1/3	OK: eth1: Link is up at 1000Mb/s, Full duplex
	eth2	OK	05-15-2009 09:02:40	51d 20h 41m 2s	1/3	OK: eth2: Link is up at 1000Mb/s, Full duplex
	eth3	OK	05-15-2009 08:56:40	51d 20h 39m 29s	1/3	OK: eth3: Link is up at 1000Mb/s, Full duplex
	load	OK	05-15-2009 08:55:12	41d 22h 49m 19s	1/3	Load : 0.00 0.00 0.00 : OK
	lustre_health	OK	05-15-2009 09:03:54	51d 20h 36m 23s	1/3	SNMP OK - healthy
	lustre_oss	OK	05-15-2009 08:55:43	51d 20h 34m 50s	1/3	SNMP OK - 4
	ntp	OK	05-15-2009 08:56:24	50d 17h 5m 20s	1/3	NTP OK: Offset 0.004905 secs
	ping	OK	05-15-2009 09:01:14	9d 19h 26m 21s	1/3	PING OK: Packet loss = 0% RTT = 0.14 ms
	ssh	OK	05-15-2009 08:56:25	51d 20h 36m 23s	1/3	SSH OK - OpenSSH_4.3p2-6.cern-hpn-CERN-4.3p2-6.cern (protocol 1.99)
	swap	OK	05-15-2009 08:56:22	23d 18h 5m 20s	1/3	Swap OK: 0.00% used
	ups	OK	05-15-2009 09:01:03	51d 20h 36m 23s	1/3	UPS OK: Voltage = 230.0V Frequency = 50.00Hz Load = 0.00%

Lustre SNMP

State Breakdowns For Host Services:

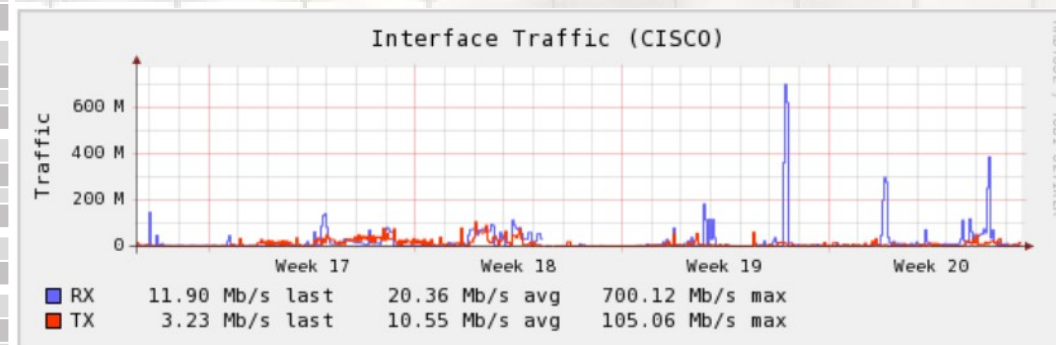
Service	% Time OK	% Time Warning	% Time Unknown	% Time Critical	% Time Undetermined
backup	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
load	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
lustre_health	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
lustre_mdt	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
lustre_mdt_osc	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
ntp	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
ping	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
ssh	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
swap	99.966% (99.966%)	0.000% (0.000%)	0.034% (0.034%)	0.000% (0.000%)	0.000%
ups	100.000% (100.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000% (0.000%)	0.000%
Average	99.997% (99.997%)	0.000% (0.000%)	0.003% (0.003%)	0.000% (0.000%)	0.000%

mdt01 31 days

Nagios@Coimbra

Service Status Details For All Hosts

Host	Service	Status	Last Check	Duration	Attempt	Status Information
CISCO	PING	OK	05-17-2009 12:32:27	104d 18h 50m 4s	1/2	OK - 192.168.4.1: rta 0.165ms, lost 0%
	TRAF-LAN	OK	05-17-2009 12:36:00	0d 16h 7m 28s	1/2	OK: GigabitEthernet0/2 is UP at 1Gbps. RX=604.2Kbps (0.06%), TX=41.13Kbps (0%)
	TRAF-RCTS	OK	05-17-2009 12:36:00	0d 16h 7m 29s	1/2	OK: GigabitEthernet0/1 is UP at 1Gbps. RX=41.04Kbps (0%), TX=604.4Kbps (0.06%)
ENV1	HUM	OK	05-17-2009 12:32:30	52d 20h 23m 6s	1/2	SNMP WARNING - *30* 35
	PING	OK	05-17-2009 12:36:08	0d 12h 30m 20s	1/2	OK - 192.168.3.3: rta 3.919ms, lost 0%
	SSH	OK	05-17-2009 12:35:46	53d 0h 36m 24s	1/2	SSH OK - cryptlib (protocol 2.0)
	TEMP	OK	05-17-2009 12:32:31	52d 19h 57m 32s	1/2	SNMP OK - 13 26 24
ENV2	HUM	OK	05-17-2009 12:34:10	38d 17h 19m 1s	1/2	SNMP OK - 55 37
	PING	OK	05-17-2009 12:31:48	12d 2h 54m 40s	1/2	OK - 192.168.3.4: rta 4.102ms, lost 0%
	SSH	OK	05-17-2009 12:32:33	38d 21h 14m 19s	1/2	SSH OK - cryptlib (protocol 2.0)
	TEMP	OK	05-17-2009 12:34:12	38d 17h 17m 22s	1/2	SNMP OK - 12 18 23
FORCE10	PING	OK	05-17-2009 12:35:50	104d 18h 50m 2s	1/2	OK - 192.168.4.2: rta 0.497ms, lost 0%
	TRAF	OK	05-17-2009 12:36:00	0d 16h 6m 59s	1/2	OK: GigabitEthernet 0/47 is UP at 1Gbps. RX=47.1Kbps (0%), TX=637.7Kbps (0.06%)
GRID001	NTP	OK	05-17-2009 12:34:14	33d 17h 32m 33s	1/2	NTP OK: Offset 7.6e-05 secs
	PING	OK	05-17-2009 12:35:52	108d 10h 51m 7s	1/2	OK - 192.168.2.1: rta 0.140ms, lost 0%
	SSH	OK	05-17-2009 12:32:37	108d 10h 4m 13s	1/2	SSH OK - OpenSSH_4.3p2-6.cern-hpn-CERN-4.3p2-6.cern (protocol 1.99)
GRID002	DNS	OK	05-17-2009 12:35:16	5d 19h 1m 12s	1/2	DNS OK: 0.059 seconds response time. www.google.com returns 74.125.79.104,74.125.79.147,74.125.79.99,74.125.79.103
	PING	OK	05-17-2009 12:35:54	106d 16h 4m 17s	1/2	OK - 192.168.2.2: rta 0.144ms, lost 0%
	SSH	OK	05-17-2009 12:32:40	108d 10h 6m 40s	1/2	SSH OK - OpenSSH_4.3p2-6.cern-hpn-CERN-4.3p2-6.cern (protocol 1.99)



17.04.09 12:38 - 17.05.09 12:38
25 May 2009, Umea Sweden

Ganglia@Lisbon

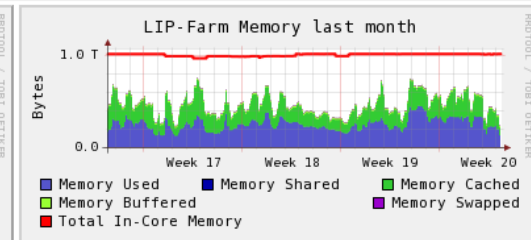
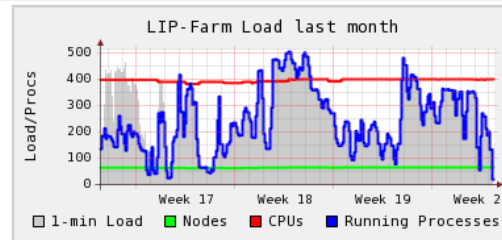
- Ganglia clusters:

- “Lustre” OSS servers

LIP-Farm (physical view)

CPU's Total: **400**
 Hosts up: **65**
 Hosts down: **0**

Avg Load (15, 5, 1m):
 1%, 18%, 17%
 Localtime:
 2009-05-15 10:25

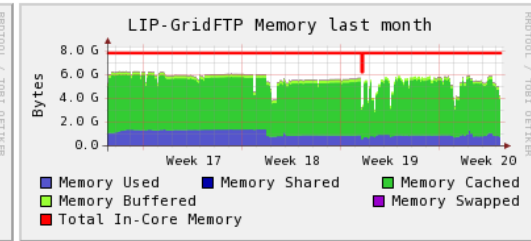
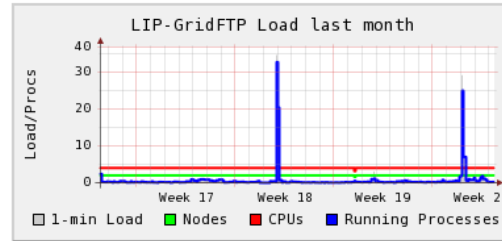


- “SRM”

LIP-GridFTP (physical view)

CPU's Total: **4**
 Hosts up: **2**
 Hosts down: **0**

Avg Load (15, 5, 1m):
 1%, 3%, 10%
 Localtime:
 2009-05-15 10:25



- “GridFTP” servers

- “Farm” - WN's

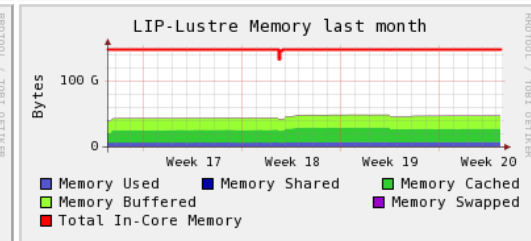
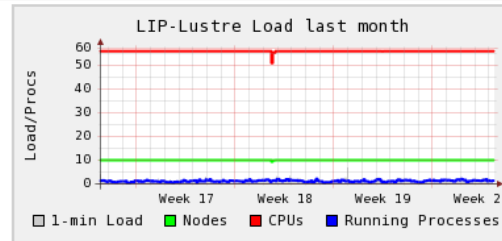
- “SGE” master

- “CE”

LIP-Lustre (physical view)

CPU's Total: **56**
 Hosts up: **10**
 Hosts down: **0**

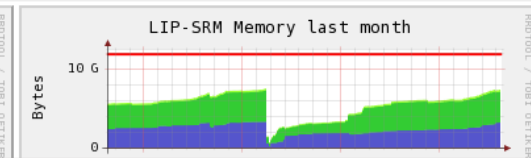
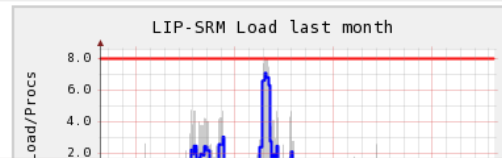
Avg Load (15, 5, 1m):
 6%, 9%, 8%
 Localtime:
 2009-05-15 10:25



LIP-SRM (physical view)

CPU's Total: **8**
 Hosts up: **1**
 Hosts down: **0**

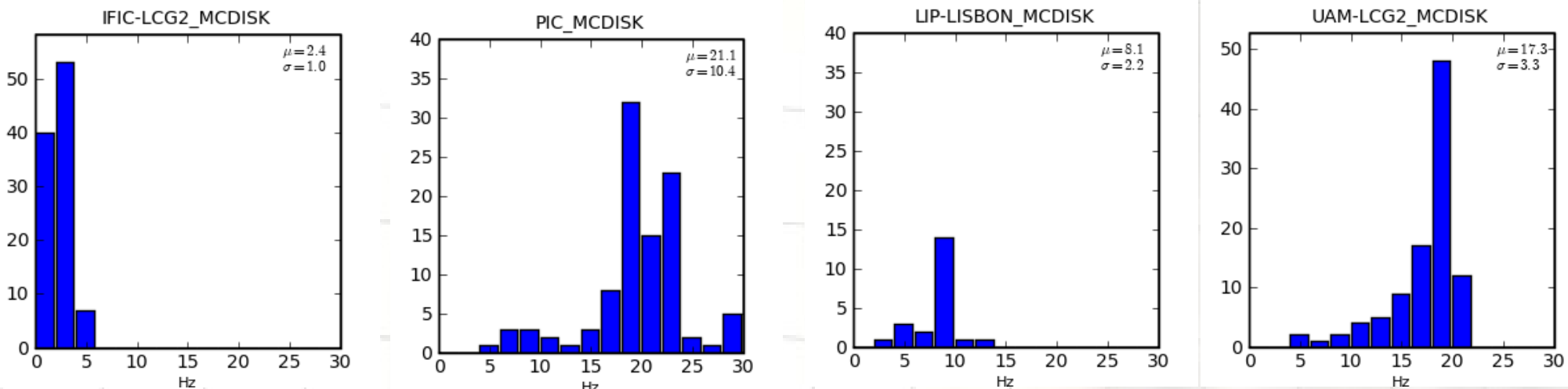
Avg Load (15, 5, 1m):
 3%, 4%, 3%



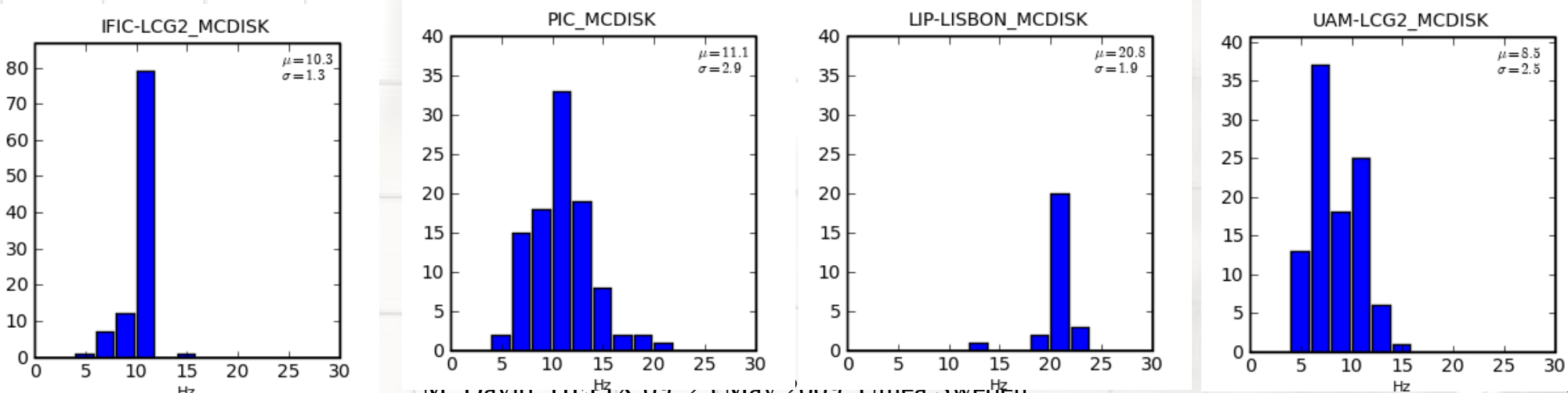
ATLAS Distributed analysis tests: HammerCloud

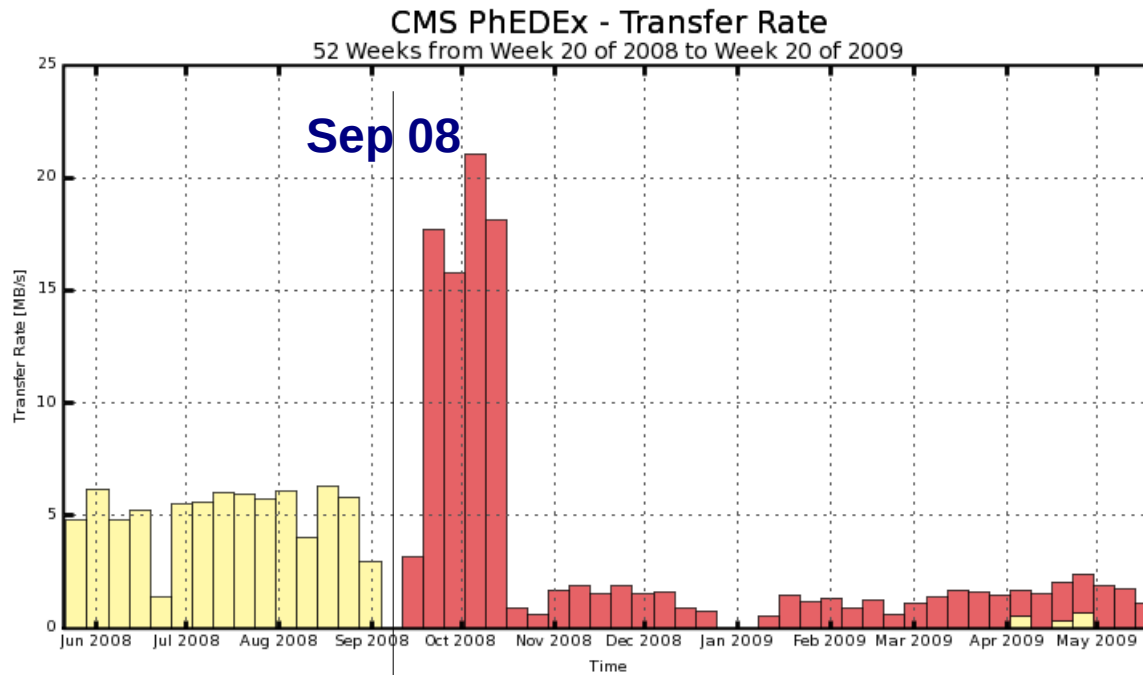
Events /second

Input method: FileStager



Input method: DQ2 Local

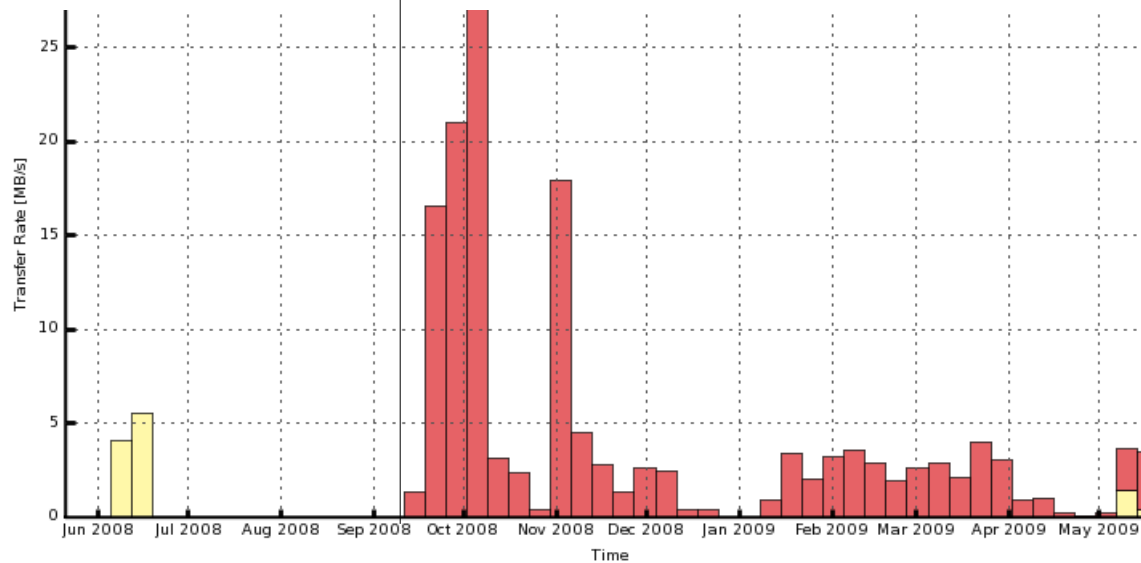




■ T2_PT_LIP_Lisbon **Coimbra**

Lisbon

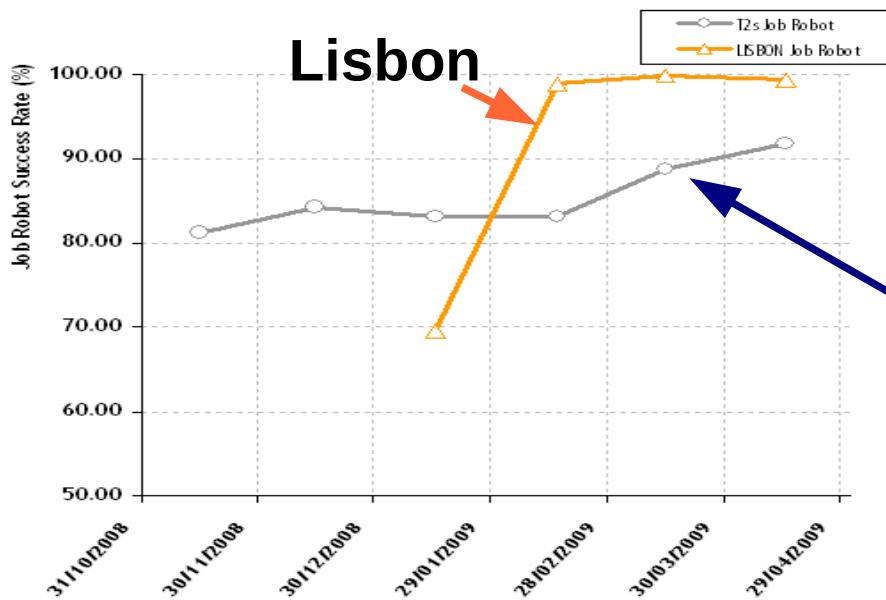
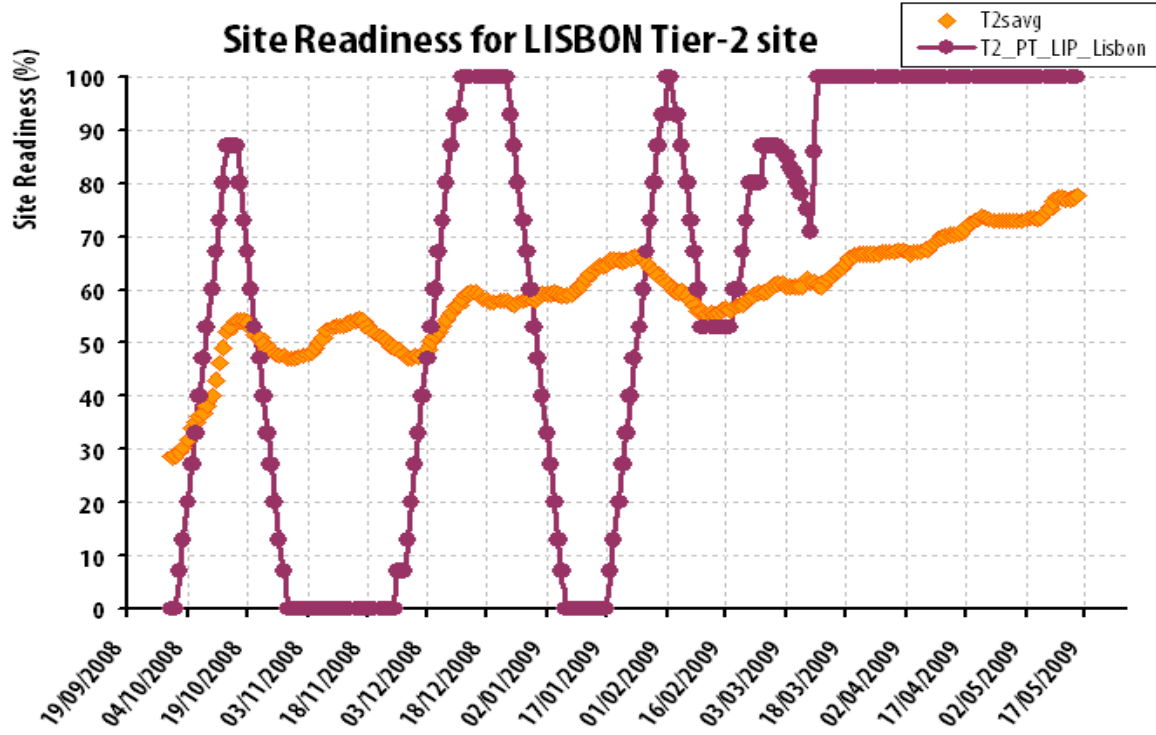
Maximum: 21.05 MB/s, Minimum: 0.00 MB/s, Average: 3.63 MB/s, Current: 1.08 MB/s



■ T2_PT_LIP_Lisbon

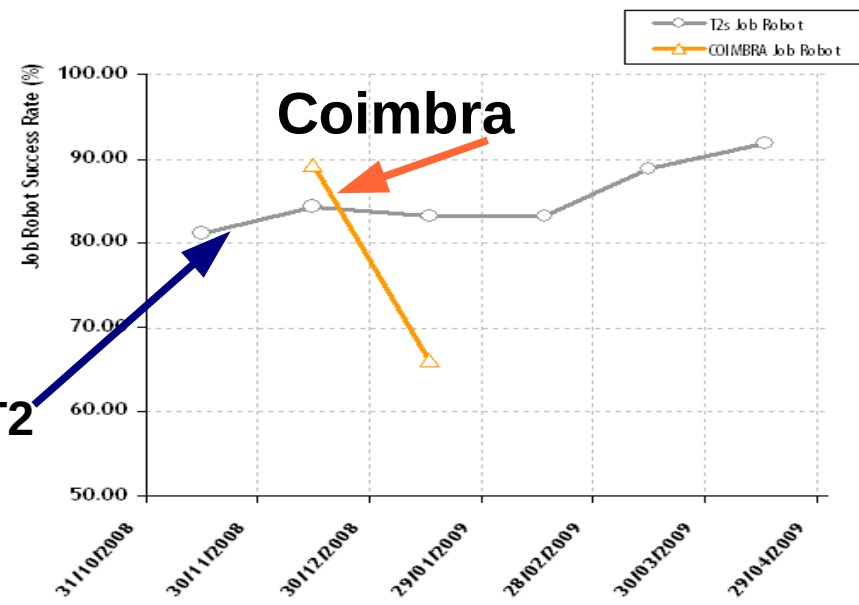
■ T2_PT_LIP_Coimbra

Maximum: 27.77 MB/s, Minimum: 0.00 MB/s, Average: 2.94 MB/s, Current: 3.48 MB/s



Job robot effic.

Avg. all T2



- The Storage Element based on the Lustre FS and StoRM SRM:
 - **Has proved to be stable, reliable and robust.**
 - **Appropriate to any site size (no nearline).**
 - **Number of updates is low:**
 - **We have made only one update of Lustre and one of StoRM since we have the system in production.**

- New hardware is arriving soon at LIP, includes:
 - **10Gb NICS for all the Lustre OSS's.**
 - **Machines with 10Gb NIC's to serve as GridFTP servers for Lisbon and Coimbra.**
- New site is now being deployed as part of the federated Portuguese T2.
- Test cluster with Lustre 1.8.0
- Development/Inclusion of other Nagios sensors and Ganglia metrics.
- STEP'09 next week will stress all we have.
 - **Note: STEP09 computing stress testing of all LHC experiments as near the real LHC data taking and processing, as possible.**

A man with a shaved head, wearing a light blue t-shirt and blue jeans, stands in a river. He has his hands on his hips and is smiling. The river is wide and calm, with a concrete bridge spanning across it in the background. The left bank is rocky and has some trees. The sky is bright and clear.

Questions??