



# T2/T3 feedback on setting-up and operating storage for the Alice VO

Jean-Michel BARBET

Laboratoire SUBATECH IN2P3/Ecole des Mines/Université de Nantes  
Nantes, France



# Outline

- Introduction
- Storage Hardware
- Network
- Xrootd deployment
- Monitoring
- Availability of storage
- Conclusion



# Introduction

- This presentation is based on feedback provided by T2/T3 site administrators answering a list of general questions
- Thank you to all site admins who took time to answer and send their comments

# ALICE jobs in the 1<sup>st</sup> year of LHC data taking at T2s

- In accordance with Alice Computing Model, T2 have to run Monte-Carlo Production and End User Analysis
- How large are T2 resources ?

		2010 (RRB year)					
		T0	CAF	T1	T2	SUM	T2/SUM %
<b>CPU</b> <b>(kHEP06)</b>	Requested	40,1	13,7	44,6	69,5	167,9	41.4
	Pledged	36,3	10,5	45,6	71,4	163.8	43.9
<b>DISK</b> <b>(TB)</b>	Requested	2313	162	5181	4879	12535	38.9
	Pledged	5500	340	6212	6089	18141	33.6

This Slide is from Galina Shabratova



# The questions were :

- how do you deal with storage installation/provisioning, difficulties ?
- how do you perceive the performances of the storage, how did you measure it ?
- how to monitor usage and failures ?
- solutions to strengthen the storage and avoid disruptions ?
- future plans
- all kind of problems and issues related to storage for the Alice VO

# Storage Hardware and filesystems

- Nothing really special here :
  - Mostly DAS boxes (10-24TB)
  - SAN solutions (at least 1, maybe more)
  - Use of RAID5 or RAID6 if SATA disks
  - Big RAID to be partitioned
- Filesystems
  - Ext3, others ?
  - Lustre at GSI (1PB local storage) and INFN



# Network

- Several sites are wondering about what is a balanced network
  - The majority of the traffic seen by the storage is from/to the workers of the same site
  - Where are the bottlenecks ?
  - Do we have recommendations ?

# Network : an attempt for guidelines

- Eygene : bw of 5MB/s per core to storage as a rule of thumb, Costin (based on ML measurements ) found about the same
- With a farm of 100 core, this gives 500Mbyte (5Gbits/s)
- Divide the total bw by the number of xrootd servers to get the minimum bw per server ? Experience show 1Gb/s/server is not enough but 10Gb/s is not needed => link aggregation
- Avoid filtering/routing between worker nodes and storage
- Is 1Gbit/s OK for an 8 core worker ? What if we go to 16 cores ?





# Xrootd

- A mix of xrootd-only storages and DPM enabled xrootd (half a dozen of sites)
  - DPM/xrootd was the solution for small sites
  - Later Alice pushed xrootd-only
  - Sites supporting several VO would prefer to have only one solution for all Vos
    - To have less different services to maintain
    - To spread traffic on more servers
  - Many sites moved nevertheless from DPM/xrootd to xrootd-only
  - DPM development team found responsive but still some problems (described later)



# Xrootd Deployment

- DPM/xrootd can be installed by YUM+YAIM (as other gLite services) or using Quattor with the QWG templates
- Xrootd-only can be installed with automatic installer in userland or with Quattor (with RPM and using the root account in this case)
- In general site admins find the installation of xrootd-only very easy



# Xrootd Maintenance

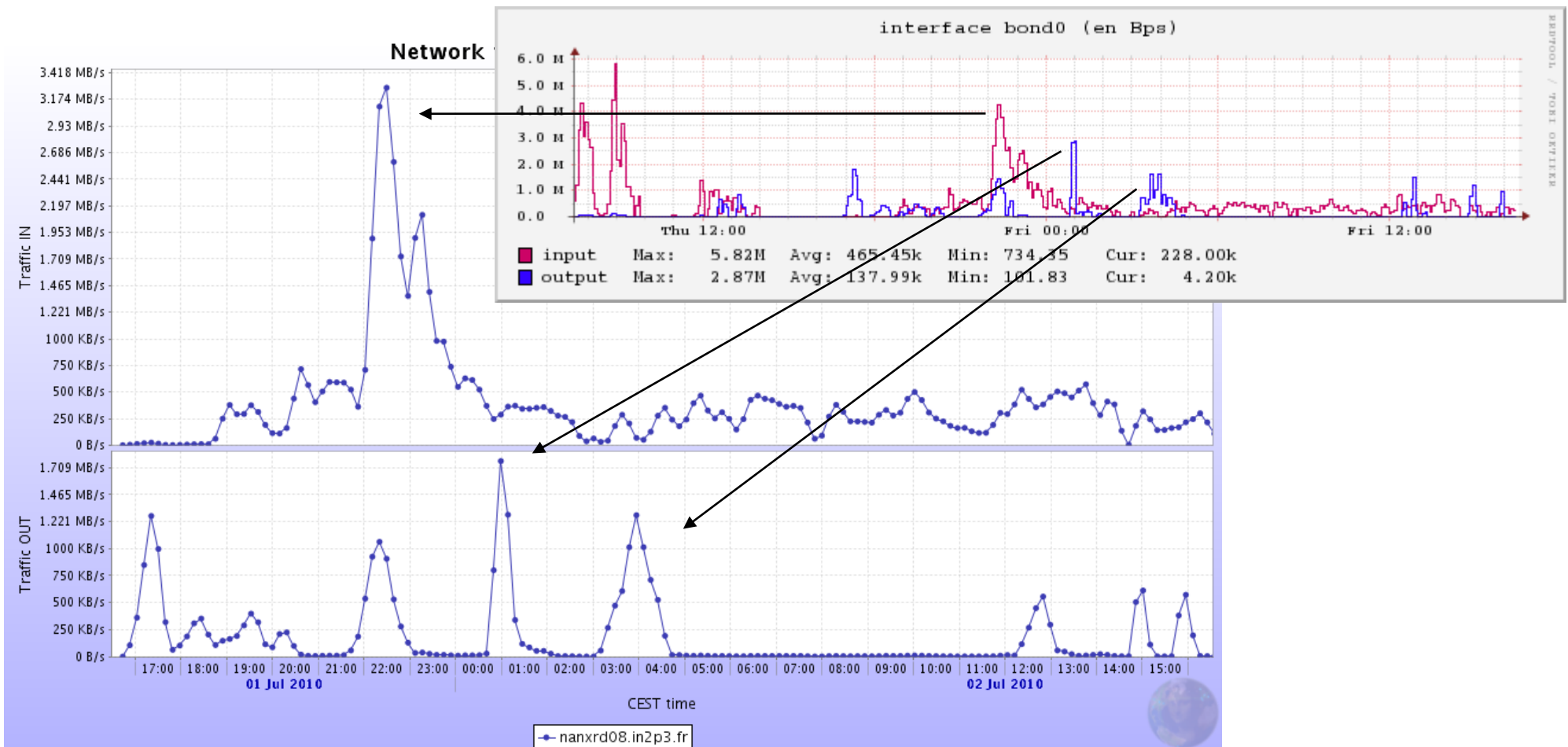
- Updates are quite easy
  - Especially if installed with installer but beware to be able to backtrack if necessary
  - RPM based installation need a new RPM to be built. The xrd-rpmer script is available for that
- A scheme for updates ?
  - Update one server, see if it still work
  - Update the manager
  - If OK, update all other servers



# Monitoring

- Monitoring can help detect failures but also shows performance
- Both Monalisa and local monitoring are used on some sites
- Local monitoring
  - Using Nagios with standard and custom probes
  - Nagios graphs (disk usage, cpu, network usage)
  - Zabbix
  - Tools that query the network fabric (SNMP, MRTG, or manufacturer specific)

# Comparing MonaLisa and local monitoring

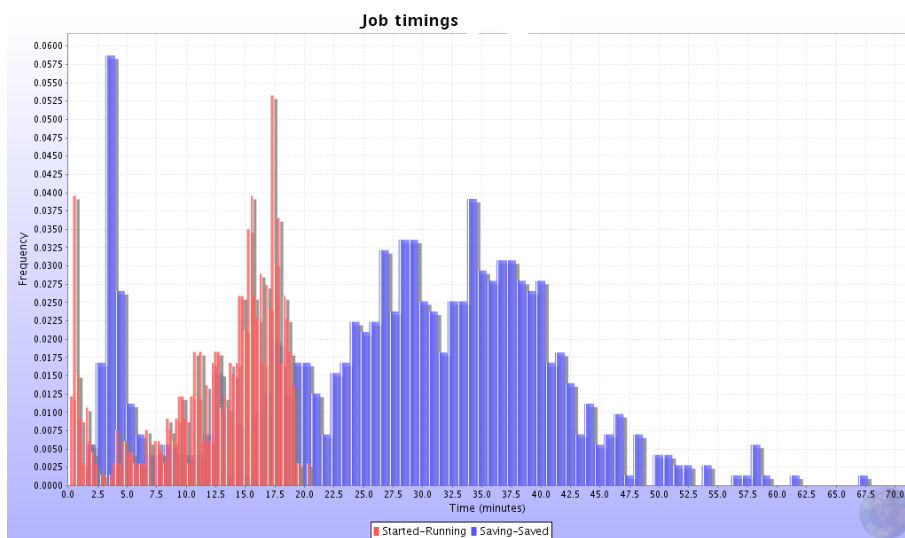


# Availability of the storage

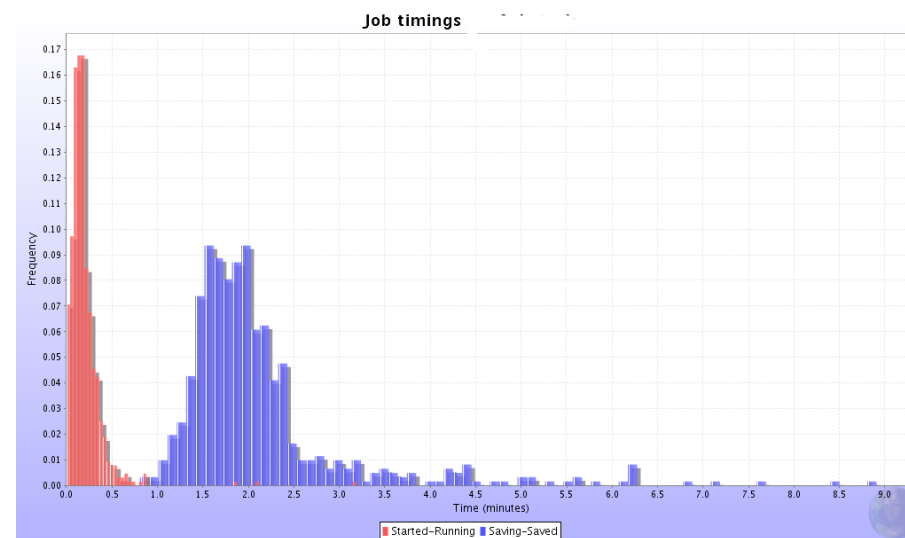
- The storage have to be available and working
  - Even more important than computing
  - Most sites are using techniques to prevent or mitigate hardware failures, possible solutions are :
    - RAID5 or RAID6
    - Redundant power adapters
    - Bonding of 2 or more Ethernet adapters
  - As a key component, the xrootd manager might be also duplicated
  - Of course good infrastructure (power, cooling) is of prior importance

# Can we visually see good storage ?

## Jobs timing profiles from MonaLisa Monitoring (\*)



A not so good site : reading and writing to the storage takes time



A good site : reading and writing to the storage is more concentrated in time

This has to be explored further, especially dependency on the workload

(\*) Credit to Costin Grigoras for pointing this



# Problems/Issues reported 1

- It seems our DPM is written into but almost never read...
- Xrootd daemon have to be restarted on DPM/xrootd
- Alice only sites have install DPM (or an SRM) only to pass the tests (2 sites)
- Xrdcp a file with a long filename fails on some Ses, seems it happens with DPM SEs
- The xrootd storage is not accounted by WLCG's current tools



# Problems/Issues reported 2

- A site without storage run no jobs
- How to drain a xrootd server in order to replace it ?
- Data in the local SE cannot be deleted in a proper way using AliEn
- Filenames in the SE are cryptic which makes it difficult to put them in correspondance with those in the local storage (GSI)
- Case of WN on a NATted private network and storage on a public network



# Conclusion (Sites wish list ?)

- Understanding the Alice activities and the workload
- Guidelines to balance worker nodes and storage and choose the appropriate network setup
- Measuring the site's performance, what to look at ?
- Distinguishing important messages to sites from ordinary talks in the mailing list