# **INFN-T1** site report

Andrea Chierici, Vladimir Sapunenko On behalf of INFN-T1 staff HEPiX spring 2012

# Outline

- Facilities
- Network
- Farming
- Grid and Middleware
- Storage
- User experience

# **INFN-Tier1: numbers**

- 1000 m<sup>2</sup> room with capability for more than 120 racks and several tape libraries
  - □ 5 MVA electrical power
  - □ Redundant facility to provide 24hx7d availability
- Within May (thanks to new 2012 tenders) even more resources
  - 1300 servers with more than 10000 cores available
  - □ **11 PBytes** of disk space and **14 PBytes** on tapes
    - Aggregate bandwith to storage: ~ 50 GB/s
- WAN link at 30 Gbit/s
  - 2x10 Gbit/s over OPN
  - □ With forthcoming GARR-X bandwidth increase is expected
- > 20 supported experiments
- ~ 20 FTE (currently evaluating some positions)



## Facilities









# Network







# Farming



# **Computing resources**

#### Currently 120K HS-06 □~9000 job slots New tender will add 31,5K HS-06 □ ~4000 potential new job slots □ 41 enclosures, 192 HS-06 per mobo We host other sites □ T2 LHCb T3 UniBO



# New tender machine

#### 2U Supermicro Twin square

- □ Chassis: 827H-R1400B
  - (1+1) redundant power supply
  - 12x 3.5" hot-swap SAS/SATA drive trays (3 for each node)
  - Hot-swappable motherboard module
- Mobo: H8DGT-HF
  - Dual AMD Opteron<sup>™</sup> 6000 series processors
  - AMD SR5670 + SP5100 Chipset
  - Dual-Port Gigabit Ethernet
  - 6x SATA2 3.0 Gbps Ports via AMD SP5100 controller, RAID 0, 1, 10
  - 1x PCI-e 2.0 x16
- AMD CPUs (Opteron 6238) 12 cores, 2,6Ghz
- 2 2tb sata hard disks 3,5"
- 40GB Ram
- 2 1620w 80gold power supply











# Issues for the future

- We would like to discuss problems with many-cores architectures
  - □1 job per core
  - □ RAM per core
  - □ Single gigabit ethernet per box
  - Local disk I/O
  - □ Green Computing: how to evaluate during tender procedures



# Grid and Middleware

# INFN

# Middleware status

# Deployed several EMI nodes Uls, CreamCEs, Argus, BDII, FTS, Storm, WNs

- Legacy glite-3.1 phased-out almost completely
- Planning to completely migrate glite-3.2 nodes to EMI within end of summer
- Atlas and LHCb switched to cvmfs for software area

□ Facing small problems with both

Tests ongoing on cvmfs server for SuperB



# WNoDeS: current status

- WNoDeS is deployed in production for two VOs
  - □ Alice: no need for direct access to local storage
  - Auger: needs a customized environment requiring a direct access to a mysql server.
- Current version of WNoDeS deployed in production is using the GPFS/NFS gateway.



# WNoDeS: development

- WNoDeS will be distributed with EMI2
- New feature called mixed mode
  - □ Mixed mode avoids to statically allocate resources to WNoDeS
  - A job can be executed on a virtual resource dynamically instantiated by WNoDeS
  - The same resources (the same hv) can be used to execute standard jobs
  - □ No resources overbooking
  - No hard resource limit enforcement which will be provided using cgroup
- General improvements

# Docet (Data Operation Center Tool)

- DB-based webtool designed and implemented internally
- In use at various INFN sites.
  - PostgreSQL, tomcat+java
  - Cmdline tool (rely on XMLRPC python webservice)
- Provides inventory for HW, SW, Actions & docs
- Initially populated by grabbing and cross-relating info from heterogeneous authoritative sources (Quattor, DNS, DHCP, XLS sheets, plaintext files) through a bunch of custom python scripts

# Docet (T1 tools)

## Cmdline Tools:

- dhcpd.conf handling: preserves consistency (update from/to docet DB to be activated)
- Cmdline query/update tool. It provides:
  - re-configuration data for nagios
  - add/list/update HW failures
  - Support for batch operations (insert new HW, remove dismissed HW)
- We plan to deploy more docet-based configuration tools

# Docet (T1 tools)

Admin Locations Components Configurations Administration Devices Network

	101-07	101-09	101-10	101-11	101	13 <sup>-t2-new.pn</sup>	ig – KSnaj <b>101</b> 014	
	42	42	42	42	4	2	42	
	sw-101-07	Tree View						
	40	II CC VICW	ann fach i Colain an Colain an Colain an Colain	Tand Could Indee				
	39	tsm 🔻 🏢 INEN."	Device					
	tsm-hsm-7	DEFAULT fakeroom Sala1						
	tsm-hsm-6							
	diskserv-san-95				Name	gridftp-	Tier3-server	
	diskserv-san-94				Hostame	ds-05-	01	
	34				Position	1		
	33		03-08		Hard Configuration	uration DELL PowerEdge M		
	nagios-storage	► <b>I</b>	▶ 08-02			Shelf 101-02		
	nagios-storage-2	storm						
	tsm-hsm-11	taj P 100-01			Edit Delete			
	tsm-hsm-10							
	diskserv-san-87							
	diskserv-san-86		101-01			Network Connections		
	diskserv-san-80	· · · · · · · · · · · · · · · · · · ·	101-02		HOSTNAME	MA	с	
	diskserv-san-64	rtagio 🕨 🕨	ds-05-01 - gridftp-Tier3-serv	er	ds-05-01	00:	1E:C9:EB:21:B8	
	storm-he-atlas-03	Storm-			COR Lake			
	storm-be-lbcb-03	aridi	ds-05-02 - gridftp-Tier3-serv	<	Capture m	nde Rectano		
	diskserv-san-80	storm 🕨 🗐	ds-05-03 - apfs-Tier3-server					
Ē	tsm-hsm-13		service spice record donter		Groups Associated	Snapshot <u>d</u>	elay: No delay	
	diskserv-san-78		ds-05-04 - gpfs-Tier3-server		NAME			



# Storage

# Storage resources

- 8.4 PB of on-line disk with GEMSS
  - □ 7 **DDN** S2A 9950
    - 2 TB SATA for data, 300 GB SAS for metadata
  - □ 7 **EMC<sup>2</sup>** CX3-80 + 1 **EMC<sup>2</sup>** CX4-960 (1 TB disks)
  - 2012 acquisition: 3 Fujitsu Eternus DX400 S2 (3 TB SATA)
- Servers
  - ~32 NSD servers (**10 Gbps** ethernet) on DDN
  - □ ~60 NSD servers (1 Gbps ethernet) on EMC<sup>2</sup>
- Tape library SI8500 (14 PB on line) with 20 T10Kb drives and 10 T10Kc drives
  - □ 9000 x 1 TB tape capacity, 1 Gbps of bandwidth for each drive
  - □ 1000 x 5 TB tape capacity, 2 Gbps of bandwidth for each drive
  - Drives interconnected to library and tsm-hsm servers via dedicated SAN (TAN)
  - □ TSM server common to all GEMSS instances
- All storage systems and disk-servers are on SAN (FC4/FC8)

INFA



# GEMSS: Grid Enabled Mass Storage System

- Integration of GPFS, TSM and StoRM
- Our choice is driven by need to minimize management effort:
  - □ Very positive experience for scalability so far;
  - Large GPFS installation in production at CNAF since 2005 with increasing disk space and number of users;
- Over 8 PB of net disk space partitioned in several GPFS clusters served by less than 100 disk-servers (NSD + gridFTP);
  - □ 2 FTE employed to manage the full system;
  - All experiments at CNAF (LHC and non-LHC) agreed to use GEMSS as HSM

# **GEMSS** evolution

- New component in GEMSS: DMAPI Server
  - Used to intercept READ events via GPFS DMAPI and re-order recalls according to the files position on tape;
  - "Preload library" is not needed anymore;
  - Available with GPFS v 3 x



Disk-centric system with five building blocks

- 1. GPFS: disk-storage software infrastructure
- TSM: tape management system 2.
- StoRM: SRM service 3
- TSM-GPFS interface 4

Meta-data flow

Data flow

Globus GridFTP: WAN data transfers 5

	IBM components INFN components

# **GEMSS:** Timeline



### GEMSS is now used by all LHC and non-LHC experiments in production for all Storage Classes

24-apr-2012

Andrea Chierici



# User experience



# Resource usage per VO



# Jobs

![](_page_29_Figure_1.jpeg)

localhost - GridJobs Collection Stats

# LHCb feedback

- More than 1 million jobs executed
  Analysis: 600k, Simulation 400k
  - User analysis not very efficient (about 50%): too large bandwidth requested
    - available bandwidth for LHCb will be significantly raised with 2012 pledges
- Stability of the services during last year
  - □ Small fraction of failed jobs
  - Good performance of data access, both from tape and disk

INFI

# CMS feedback

- 5 Kjobs/day average (50 Kjobs/day peak)
- Up to 200 TB data transfers per week both in and out
- Recent tests proved the possibility to work with 10 GB files in a sustained way

INFI

# Alice feedback

- 6.93×10<sup>6</sup> KSI2k hours consumed in one year of very stable running
- The cloud infrastructure based on WNoDes provided excellent flexibility in catering to temporary requirements (e.g. large memory queues for PbPb event reconstruction).
- CNAF holds 20% of ALICE RAW data (530 TB on tape)
- All data are accessed through an XROOTD interface over the GPFS+TSM underlying file system

INFI

![](_page_33_Picture_0.jpeg)

## Questions?