



# LHCb Report

from **09/2012** to **11/2012**

Vedaee Aresh (15/11/2012)

# Summary



- **General Status and Updates**
- **LHCb and Cloud**
- **LHCb pilots issues**

# General status and updates (1/2)



- **CVMFS (18/09/2012):**
  - sw cache 10GB
  - cache-conddb 1.5GB
- **HPSS (08/10/2012):** additional 400TB of tape
- **LFC Server (17/10/2012):** lfc-lhcb-ro.in2p3.fr decommission
- **DISK (before 09/2012):** enabled checksum (based on gfal library) for FTS and file upload



# General status and updates (2/2)



## ■ Computing

- 2300 running jobs in avg
- ~4500 ended jobs per day (last day 8k)
- ~50% pilots wallclock\_time  $\leq$  1000s
- ~20% pilots wallclock\_time  $\geq$  30h (huge)

## ■ Disk

- 12% free space on LHCb-Disk

## ■ Tape

- 250TB consumed since 10/2012
- 150TB left out of pledged 1400TB

# ► LHCb and Cloud (1/2)



## ■ Why the cloud?

VM tailored for LHCb's needs to monitor/investigate critical jobs in a isolated environment. E.g. LHCb jobs and memory consumption.

## ■ with Batch System

- Dedicated VM cluster (mem, scratch, swap, etc.) and queue
- Ensure safe interact. with external services (dcache, cvmfs, etc.)
- Test and monitor (see Dirac Web Portal)

## ■ without Batch System

- **Implementation** (<http://dl.acm.org/citation.cfm?id=2116173>)

- 1.start manually the vm and run dirac pilots
- 2.spawn automatically the vm

- **Current status**

- Test machine: cctbdirac01 + SL6 + nova (openstack api) client
- Cloud VM: SL6 + cvmfs + dirac client

- **TODO**

- Testing dirac and lhcb pilots: network access still missing (open bu problems with openstack config)
- OCCl development ready in 1 or 2 months

# LHCb Pilots Conditions (1/2)



## ■ **CONDITION at cc-IN2P3**

- Share T1/T2
  - P\_lhcb=5, P\_lhcb\_admin=1, P\_lhcb\_pilot=90, P\_lhcb\_prod=4
  - LHCb=16076, ATLAS=31067, CMS=11471
- Execution constraint: 3200 in the whole farm
- Submission constraint: 3000 jobs (queued+run) per user
  - lhcb047, lhcb049 for T1
  - lhcb097, lhcb099 for T2
- Queues and CPU/MEM limits:

Queue	CPU (soft) Limit	MEM limit
Verylong (Huge)	46h	5GB
Long	30h	4GB
Medium	5h	3GB

## ■ **LHCb Computing model**

- Tier 1: reco/stripping/merging/user (with input data)
- Tier 2: user jobs (without input data) and mc jobs

## ■ **Dirac Algorithm**

- task queue treating sites differently
- job classified according to CPU/MEM requirements
  - Reconstruction: about 36h, and 2GB Mem
  - Stripping: about 18h, and 3-4GB Mem
  - Merge: about 30min, and 4-5GB Mem
- # waiting jobs > # (submitted/waiting) pilots
- trend forecast





# Empty Pilots



- **Pilot redundancies:** efficiency vs. effectiveness
- **Match delay:** every N seconds a job of a given job type to be matched by a pilot. E.g.: if 1000 pilots, only reco jobs, reco matching delay 20sec, and pilot time out 3min, then 9 reco jobs match and the other 991 pilots die. Solution: throttle (prompt)reco jobs that are I/O consuming -> ramp-up of SE/WN transactions overloading srm servers
- **Retry delay (todo?):** randomize (+ notification)
- **Misconfiguration (fixed):** e.g. ratio waiting/running (e.g. T2 in Jul/Oct)
- **GE Bulk submission (fixed before 08/2012: 30sec scheduling passes)**
  - reducing the number of lhcb pilot submissions per GE cycle
  - increasing the frequency of the lhcb pilot submissions

## ▶ Short Pilots (1 short job)



- **Overkill:** e.g. failed due to download input sandbox or input data resolution
- **Bugs (fixed):** pilot executing a job cpu limit (normalization factor ignored due to benchmark discrepancy)
- **# pilots >> # jobs:** it should be absorbed over time



- **VMEM violation:** GE kills the whole process tree and not the application (gaudi-run.py) so no job logs
  - Additional swap on workers, not possible.
  - Likely solution: job wrapper sends back job logs when it receives the signal (SIGXCPU = 152) -> Problem: no standardized s/h signals among sites



**QUESTIONS?**