○ Reminder:

- ❏ Two-fold goal: produce and reconstruct useful data, exercise the LHCb Computing model, DIRAC and ganga

- ❏ To be tested:

  - ☆ Software distribution

  - ☆ Job submission and data upload (simulation: no input data)

  - ☆ Data export from CERN (FTS) using MC raw data (DC06-SC4)

  - ☆ Job submission with input data (reconstruction and re-reconstruction)

    - ❄ For staged and non-staged files

  - ☆ Data distribution (DSTs to Tier1s T0D1 storage)

  - ☆ Batch analysis on the Grid (data analysis and standalone SW)

  - ☆ Datasets deletion

- ❏ LHCb Grid community solution

  - ☆ DIRAC (WMS, DMS, production system)

  - ☆ ganga (for analysis jobs)

LHCB EXPERIENCE WITH SERVICES

○ **Summer 2006**

    ❑ Data production on all sites

        ✰ Background events (~100 Mevts b-inclusive and 300 Mevts minimum bias), all MC raw files uploaded to CERN

○ **Autumn 2006**

    ❑ MC raw files transfers to Tier1s, registration in the DIRAC processing database

        ✰ As part of SC4, using FTS

            ❄ Ran smoothly (when SEs were up and running, never 7 at once)
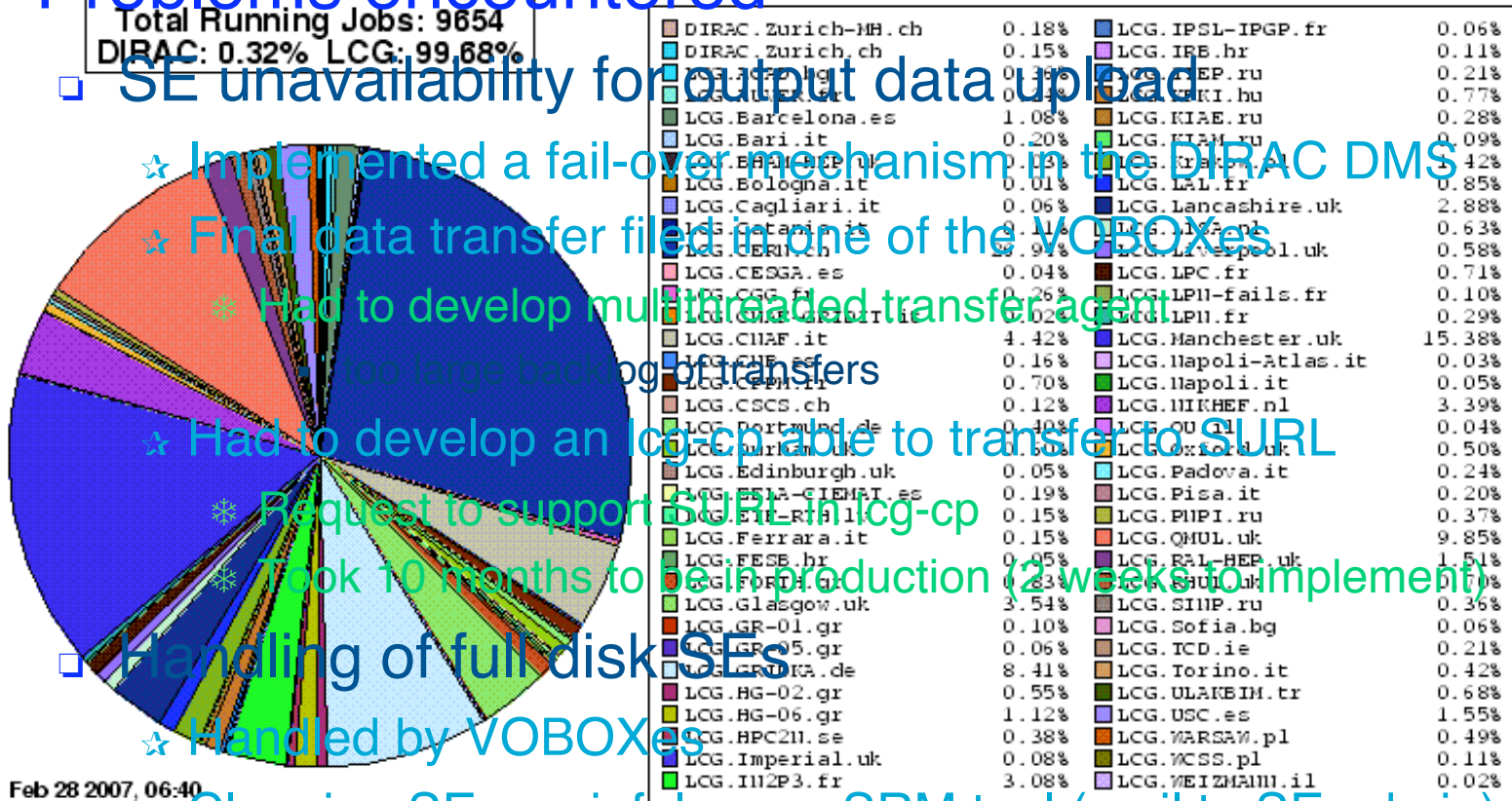
○



Averaged Throughput during the last 24 hrs (08/06 – 09/06)
Data Transfer For 'LHCb' From All Sites To All Sites

*(Legend: FZK, IN2PCC, PIC, RAL, SARA, UNREGD)*

○ **February 2007 onwards**

- ❏ Background events reconstruction at Tier1s
  - ✰ Uses 20 MC raw files as input
    - ❄ were no longer on cache, hence had to be recalled from tape
  - ✰ output rDST uploaded locally to Tier1

○ **June 2007 onwards**

- ❏ Background events stripping at Tier1s
  - ✰ Uses 2 rDST as input
  - ✰ Accesses the 40 corresponding MC raw files for full reconstruction of selected events
  - ✰ DST distributed to Tier1s
    - ❄ Originally 7 Tier1s, then CERN+2
    - ❄ need to clean up datasets from sites to free space

# Software distribution

□ **Performed by LHCb SAM jobs**

  ☆ See Joël Closier's poster at CHEP

□ **Problems encountered**

  ☆ Reliability of shared area: scalability of NFS?

  ☆ Access permissions (lhcbsgm)

  ☆ Move to pool accounts…

  ☆ Important: beware of access permissions when changing accounts mapping at sites!!!

  ❄ moving to pool accounts was a nightmare

- ○ Up to 10,000 jobs running simultaneously
  - ❑ Continuous requests from physics teams
- ○ Problems encountered
  - ❑ SE unavailability for output data upload
    - ☆ Implemented a fail-over mechanism in the DIRAC DMS
    - ☆ Final data transfer filed in one of the VOBOXes
      - ❄ Had to develop multithreaded transfer agent
        - ❄ too large backlog of transfers
    - ☆ Had to develop an lcg-cp able to transfer to SURL
      - ❄ Request to support SURL in lcg-cp
      - ❄ Took 10 months to be in production (2 weeks to implement)
  - ❑ Handling of full disk SEs
    - ☆ Handled by VOBOXes
    - ☆ Cleaning SEs: painful as no SRM tool (mail to SE admin)



Total Running Jobs: 9654
DIRAC: 0.32%  LCG: 99.68%

| | | | |
|---|---|---|---|
| DIRAC.Zurich-MH.ch | 0.18% | LCG.IPSL-IPGP.fr | 0.06% |
| DIRAC.Zurich.ch | 0.15% | LCG.IRB.hr | 0.11% |
| | | LCG.ITEP.ru | 0.21% |
| | | LCG.KFKI.hu | 0.77% |
| LCG.Barcelona.es | 1.08% | LCG.KIAE.ru | 0.28% |
| LCG.Bari.it | 0.20% | LCG.KIAM.ru | 0.09% |
| | | LCG.Krakow.pl | 0.42% |
| LCG.Bologna.it | 0.01% | LCG.LAL.fr | 0.85% |
| LCG.Cagliari.it | 0.06% | LCG.Lancashire.uk | 2.88% |
| LCG.Catania.it | | LCG.Legnaro.it | 0.63% |
| LCG.CERN.ch | 9.9% | LCG.Liverpool.uk | 0.58% |
| LCG.CESGA.es | 0.04% | LCG.LPC.fr | 0.71% |
| LCG.CGG.fr | 0.26% | LCG.LPN-fails.fr | 0.10% |
| | | LCG.LPN.fr | 0.29% |
| LCG.CNAF.it | 4.42% | LCG.Manchester.uk | 15.38% |
| LCG.CPPM.fr | 0.16% | LCG.Napoli-Atlas.it | 0.03% |
| LCG.CSCS.ch | 0.12% | LCG.Napoli.it | 0.05% |
| | | LCG.NIKHEF.nl | 3.39% |
| LCG.Dortmund.de | | LCG.OU.uk | 0.04% |
| LCG.Durham.uk | | LCG.Oxford.uk | 0.50% |
| LCG.Edinburgh.uk | 0.05% | LCG.Padova.it | 0.24% |
| LCG.ESA-CIEMAT.es | 0.19% | LCG.Pisa.it | 0.20% |
| LCG.Ferrara.it | 0.15% | LCG.PNPI.ru | 0.37% |
| LCG.FESB.hr | 0.05% | LCG.QMUL.uk | 9.85% |
| | | LCG.RAL-HEP.uk | 1.51% |
| LCG.Glasgow.uk | 3.54% | LCG.SINP.ru | 0.36% |
| LCG.GR-01.gr | 0.10% | LCG.Sofia.bg | 0.06% |
| LCG.GR-05.gr | 0.06% | LCG.TCD.ie | 0.21% |
| LCG.GRIDKA.de | 8.41% | LCG.Torino.it | 0.42% |
| LCG.HG-02.gr | 0.55% | LCG.ULAKBIM.tr | 0.68% |
| LCG.HG-06.gr | 1.12% | LCG.USC.es | 1.55% |
| LCG.HPC2N.se | 0.38% | LCG.WARSAW.pl | 0.49% |
| LCG.Imperial.uk | 0.08% | LCG.WCSS.pl | 0.11% |
| LCG.IN2P3.fr | 3.08% | LCG.WEIZMANN.il | 0.02% |

Feb 28 2007, 06:40

# Reconstruction jobs

○ Needs files to be staged

 ❑ Easy for first prompt processing, painful for reprocessing

 ❑ Developed a DIRAC stager agent

  ☆ Jobs are put in the central queue only when files are staged

○ File access problems

 ❑ Inconsistencies between SRM tURLs and root access

 ❑ problems with ROOT finding the HOME directory

  ☆ at RAL, fixed by providing an additional library (compatibility mode on SLC4)

 ❑ unreliability of rfio, problems with rootd protocol authentication on the Grid (now fixed by ROOT)

 ❑ Impossible to copy input data locally (not enough disk guaranteed)

  ☆ advise from SE experts: better access files from server…

 ❑ lcg-gt returning a tURL on dCache but not staging files

  ☆ Workaround with dccp, then fixed by dCache

# File access problems (cont'd)

○ **Some files are not retrievable from tape**

- ❑ registered in our LFC
- ❑ found using srm-get-metadata
- ❑ but fail to get a tURL (error in lcg-gt)

○ **Some files are temporarily unavailable**

- ❑ e.g. those above (in case tape is corrupted, stuck…)
- ❑ files on D1T0 that are not actually on disk
  - ☆ srm-get-metadata: isCached=false
- ❑ need to establish a protocol to get warning from site
  - ☆ will set a flag in LFC indicating the replica is temporarily unavailable (not used for matching jobs)

○ **Staging at some sites extremely slow**

- ❑ problems with SE software?
- ❑ problems of configuration?
  - ☆ number of servers, number of tape drives
- ❑ on our side, need to tune the number of stage requests issued in one go
  - ☆ try and optimise the recall from tape

# What is still missing?

- ○ gLite WMS
  - ❏ Many attempts at using it, encouraging
    - ☆ Still not used in production because of…
- ○ Full VOMS support
  - ❏ Many problems of mapping when using VOMS
    - ☆ LHCb wanted to use group/role : wasn't correctly implemented at sites
      - ❄ rolling back to "default" behavior not using groups
    - ☆ Problems of LFC registration in existing directories
      - ❄ e.g. when moving to pool accounts for production group
      - ❄ DN/FQAN changes can't be handled but by root admin
      - ❄ giving group write permission is not really optimal!
    - ☆ No castor proper authentication (i.e. no security for files)
- ○ Agreement and support for generic pilot jobs
  - ❏ Essential for good optimisation at Tier1s
    - ☆ Prioritisation of activities (simulation, reconstruction, analysis)

# Storage Resources

○ **Main problem encountered is with Disk1TapeX storage**

- ❑ 3 out of 7 Tier1s didn't provide what had been requested
  - ☆ Continuously change distribution plans for LHCb
  - ☆ Need to clean up datasets to get space (painful with SRM v1)
- ❑ Not efficient to add servers one by one
  - ☆ When all servers are full, puts a very large load on the new server
- ❑ Not easy to monitor the storage usage
  - ☆ developed a specific agent reporting every day from LFC
  - ☆ other agents checking integrity between SEs and catalogs

○ **Too many instabilities in SEs**

- ❑ Full time job checking availability
  - ☆ Enabling/disabling SEs in the DMS
  - ☆ VOBOX helps but needs guidance to avoid DoS

○ **Several plans for SE migration**

- ❑ RAL, PIC, CNAF, SARA (to NIKHEF): to be clarified

LHCb EXPERIENCE WITH SERVICES

○ LHCb happy with the proposed agreement from JSPG (EDMS 855383)

  ❏ Eager to see it endorsed by all Tier1s

    ☆ Essential as LHCb run concurrent activities at Tier1's

  ❏ DIRAC prepared for running its payload through a glexec-compatible mechanism

    ☆ Wait for sites to deploy the one they prefer

# Middleware deployment cycle

○ Problem of knowing "what runs where"

❏ Reporting problems that was fixed long ago

☆ but either were not released or not deployed

○ Attempt at getting the client MW from LCG-AA

❏ very promising solution

❏ very collaborative attitude from GD

☆ versions for all available platforms installed as soon as ready

☆ allows testing on LXPLUS and on production WNs

❊ tarball shipped with DIRAC and environment set using CMT

❊ not yet in full production mode, but very promising

☆ allows full control of versions

❊ possible to report precisely to developers

❊ no way to know which version runs by default on a WN

# SLC4 migration

- Straightforward for LHCb applications
  - problem was middleware clients used by them
    - dCache, gfal, lfc…
- Usage by DIRAC
  - binaries are OK
    - except lcg-cp that had a regression (2 weeks to find out)
  - python binding is not OK at some sites because…
- Inconsistencies between MW and OS
  - middleware is 32-bit only
  - hence WNs should by default expose a 32-bit architecture when being accessed from grid queues
    - at CERN, python is 64-bit
    - in addition unnecessary environment variables are making the case even more complicated
- DIRAC3
  - will import all necessary middleware (including python)
    - from LCG-AA, installed on sites by SAM jobs

❍ Very impractical to test client MW on PPS

  ❑ completely different setup for DIRAC

  ❑ hard to verify all use cases (e.g. file access)

❍ Was used for testing some services

  ☆ e.g. gLite WMS

  ❑ but easier to get an LHCb instance of the service

    ☆ known to the production BDII

    ☆ possibility to use or not depending on reliability

      ❊ example: slc4 CEs were needed in order to find out all pbs

    ☆ sees all production resources

      ❊ caveat: should not break e.g. production CEs

        ▪ but expected to be beyond that level of testing…

❍ PPS uses a lot of resources in GD

  ❑ worth discussing with experiments if needed…

    ☆ no definite answer to the question from LHCb…

# Monitoring & availability

❍ **Essential to test sites permanently**

  ❏ See J.Closier's poster at CHEP

  ❏ Use the SAM framework

    ✩ check availability of CEs open to LHCb

    ✩ install LHCb and LCG-AA software

      ❆ platform dependent

    ✩ reports to the SAM database

    ✩ LHCb would like to report the availability as they see it

      ❆ no point claiming a site is available just for the ops VO

  ❏ Faulty sites are "banned" from the DIRAC submission

  ❏ Faulty SEs or full disk-SEs can also be "banned" from the DMS (as source and/or destination)

# **Conclusions**

○ **LHCb using WLCG/EGEE infrastructure successfully**

    ❏ Eagerly waiting for generic pilots general scheme

○ **Still many issues to iron out (mainly DM)**

    ❏ SE reliability, scalability and availability

    ❏ Data access

    ❏ SRM v2.2

    ❏ SE migration at many sites

○ **Trying to improve certification and usage of middleware**

    ❏ LCG-AA deployment, production preview instances

○ **Plans to mainly continue regular activities**

    ❏ Move from "challenge mode" to "steady mode"