



LCG-France Tier-1 & AF

Réunion de Coordination

Fabio Hernandez fabio@in2p3.fr

Lyon, 14 mai 2009







Table des Matières



- Réunions GDB et MB de mai 2009
- Disponibilité, fiabilité, efficacité des sites
- Chantiers en cours
- Thème(s) du jour

E.Hernandez



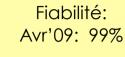
GDBs & MBs



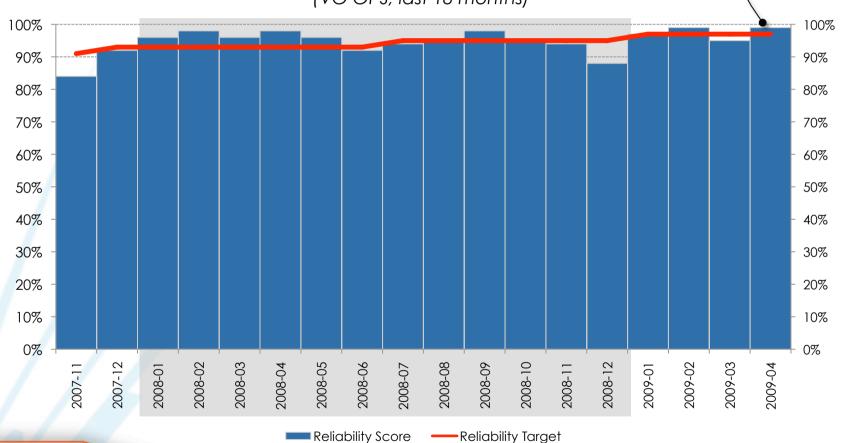
- Agendas:
 - GDB: http://indico.cern.ch/categoryDisplay.py?categld=31181
 - MB: http://indico.cern.ch/categoryDisplay.py?categld=666
- Principaux sujets traités
 - Préparation STEP'09
 - Scientific Linux 5



Fiabilité du tier-1



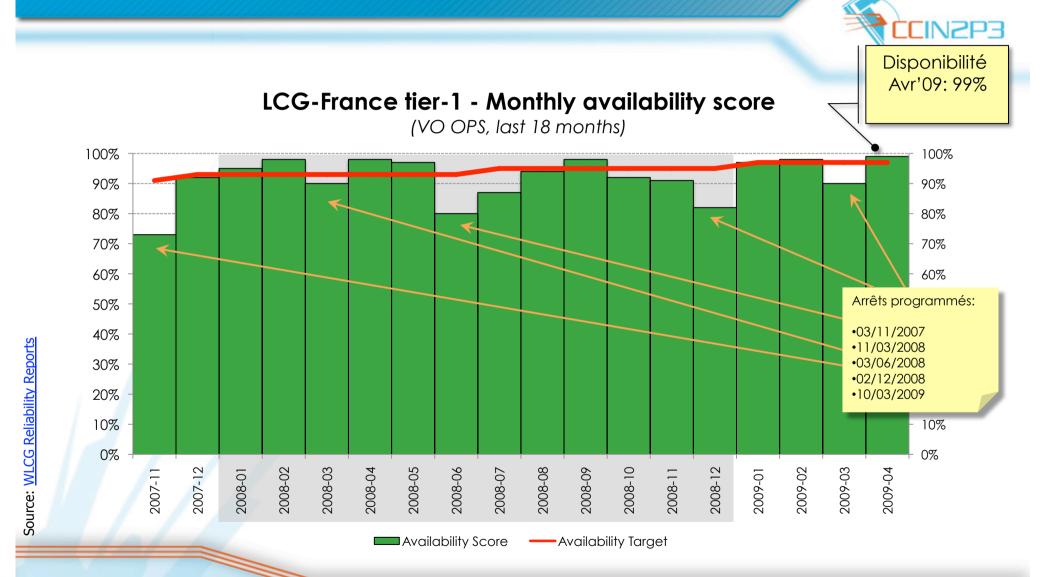




Source: WLCG Reliability Reports



Disponibilité du tier-1

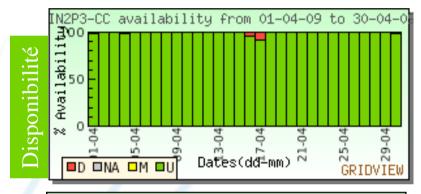


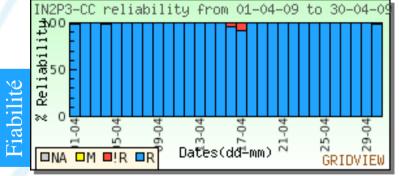


Disponibilité & Fiabilité tier-1



Avril 2009 (VO OPS)





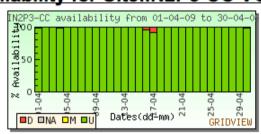


Disponibilité tier-1



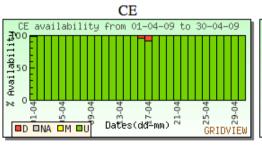
Avril 2009

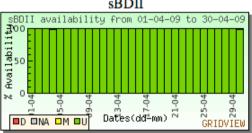
Overall Service Availability for Site: IN2P3-CC VO: OPS (Daily Report)

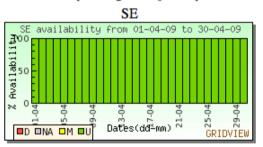


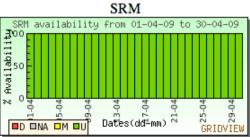
Score: 99%

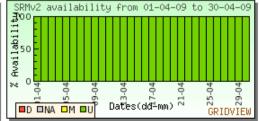
Individual Service Availability for site: IN2P3-CC VO: OPS (Daily Report)











SRMv2





Disponibilité tier-1 (suite)



Avril 2009 (suite)

Dates(dd-mm)

Service Instance Availability for site: IN2P3-CC VO: OPS (Daily Report)

Dates(dd-mm)



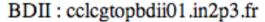


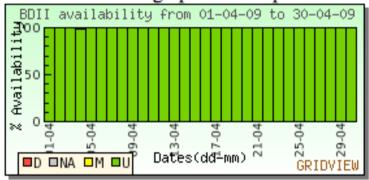
Disponibilité tier-1 (suite)

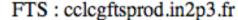


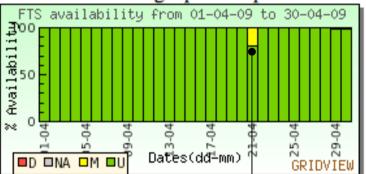
Avril 2009 (suite)

Central Services







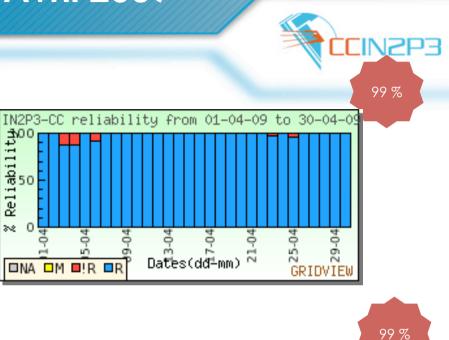


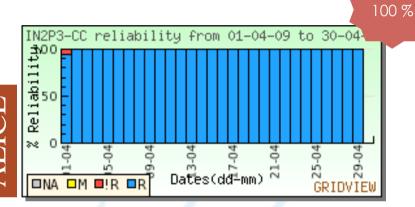
Shutdown du service pour mise en production de FTS v2.1 sous SL4 64bits

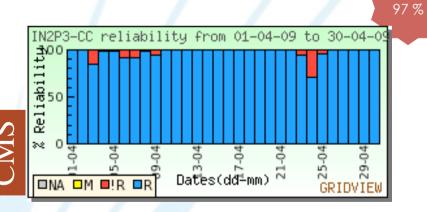
Disponibilités tier-1s + CERN

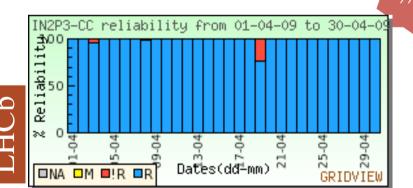


Fiabilité vue par les VOs: Avril 2009









Source: http://gridview.cern.ch





EGEE France – disponibilité & fiabilité

Ed2NISD3

Avril 09

Region	Avail- ability	Reli- ability
France	97 %	98 %
SouthWesternEurope	95 %	97 %
UKI	94 %	96 %
NorthernEurope	94 %	96 %
SouthEasternEurope	93 %	94 %
GermanySwitzerland	93 %	93 %
CentralEurope	93 %	93 %
CERN	88 %	88 %
AsiaPacific	80 %	82 %
Italy	78 %	88 %
Russia	60 %	68 %

F.Hernandez



EGEE France – disponibilité & fiabilité



Region Site
France (France)

CGG-LCG2

IBCP-GBIO IN2P3-CC

IN2P3-CC-T2 IN2P3-CPPM IN2P3-IPNL IN2P3-IRES IN2P3-LAPP IN2P3-LPC IN2P3-LPSC

IN2P3-SUBATECH IPSL-IPGP-LCG2

ESRF GRIF Ces chiffres
paraissent suspects
pour certains sites
(Phy.CPU = Log.CPU)

Phy. Log.		Avail-	Reli-	Availability History			
CPÚ		KSI2K	ability	ability	Jan-09	Feb-09	Mar-09
1							
42	42	75	98 %	98 %	98 %	100 %	98 %
80	80	49	91 %	91 %	46 %	73 %	96 %
16	16	43	86 %	86 %	83 %	89 %	98 %
3,338	2,180	3,908	100 %	100 %	99 %	100 %	92 %
10	10	5	55 %	97 %	60 %	67 %	94 %
1,074	4,296	3,832	99 %	99 %	97 %	98 %	90 %
1,074	4,296	3,832	99 %	99 %	96 %	97 %	89 %
358	358	537	98 %	99 %	99 %	97 %	95 %
452	440	656	98 %	99 %	98 %	96 %	98 %
664	628	1,526	84 %	84 %	95 %	96 %	93 %
512	512	1,133	96 %	98 %	94 %	100 %	98 %
448	448	802	84 %	99 %	94 %	93 %	99 %
120	112	43	71 %	96 %	97 %	99 %	94 %
275	380	803	97 %	97 %	99 %	96 %	99 %
34	34	41	96 %	96 %	94 %	100 %	96 %

Source: https://edms.cern.ch/document/963325

F.Hernandez



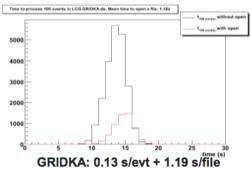
LHCb Analysis Tests

Running the Test

- Each submission consist of:
 - ~ 600 jobs, reading 100 different files with 500 events per file.
 - Average file size 200 MB
 - 1 Job = 100 Files = 50.000 events ~ 20 GB ~ 1-2 hours~ 2-3 MB/s
- InputData:
 - User submits list of LFNs
 - DIRAC checks LFC for replicas and check availability at the sites.
 - Job receives SURLs.
 - SURLs are converted to tURLs by asking to SRM (getturls)
- Application:
 - Root base application open the files in sequence, reads and analyzes all events directly from the SE (without copy to WN).

The different represents the represents the represents the representation in LCG.

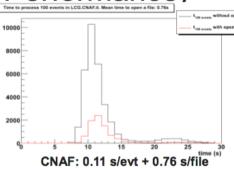
Time to process 100 events in LCG.

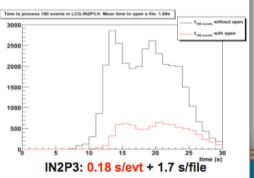


CCINSP3

Results (Detail Performance)

- Distribution of Wall Clock time to process 100 consecutive events.
- Distribution of Wall Clock time to process 100 consecutive events including a new file opening.
- The difference of the means represents the file opening time.





Source: Ricardo Graciani, GDB, 13/05/2009

-.Hernandez



LHCb Analysis Tests (suite)

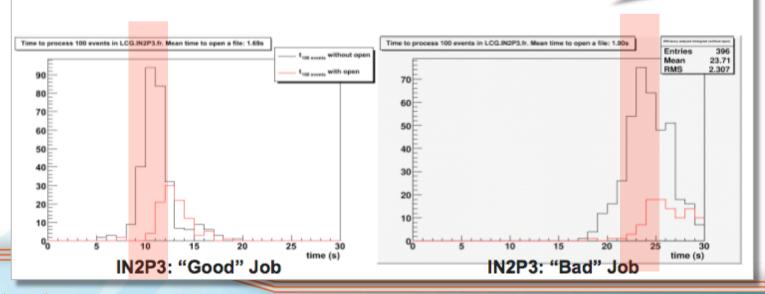


IN2P3 observations

- Poor CPU Efficiency at IN2P3 (< 50%)
- Large variation in time to process events
- Seems to be Job to Job variations
- Needs further investigations

Comportements différents des jobs. Pistes à explorer, d'après un expert dCache:

- dCache debug = 2 sur certains worker nodes
- Nombre maximum de fichiers ouverts par pool (MaxMovers?)





SL5 – Portage du code des expériences

3 - Status and Plans

ATLAS

- The goal is to have a full native build by next week (deployed on the GRID by ~August 2009)
- Will continue to produce binaries for SL4 and SL5 concurrently
- Inclined to retain the existing default (SLC4/gcc34/32-bit) as primary platform until after the 2009-2010 physics run

CMS

- Finished native port 64 bit/SL5. Everything builds, basic tests seem okay, will proceed to full validation.
- Rather switch from SL4 to SL5 binaries in one go at their convenience

LHCb

- Port not yet finalized (not yet a complete working release)
- Plan is to use native SL5/gcc43/64bit for the real data and the corresponding MC productions

ALICE

No problem. They would like to see the transition to 64 bits/SL5 to happened soon as possible

pere.mato@cern.ch May 12, 09

Etat d'avancement de compilation native sur SL5 avec compilateur gcc-4.3 (pas d'utilisation du compilateur par défaut sur SL5 qui est gcc-4.1)

Activation de SELinux a un impact sur certains librairies, y compris Oracle.

Passage en production n'est pas possible pour toutes les expériences avant septembre 2009.

Des WNs sous SL5 pour la validation restent nécessaires, avec un CE spécifique.

Source: Pere Mato, GDB, 13/05/2009



Consommation mémoire sous Linux 64bits



VMEM in 64bit Architecture

Python 64bits VSZ = 58 MB

 We observe that any process running in 64bit has a VMEM footprint of ~50MB (c.f. 5MB for a 32bit process)

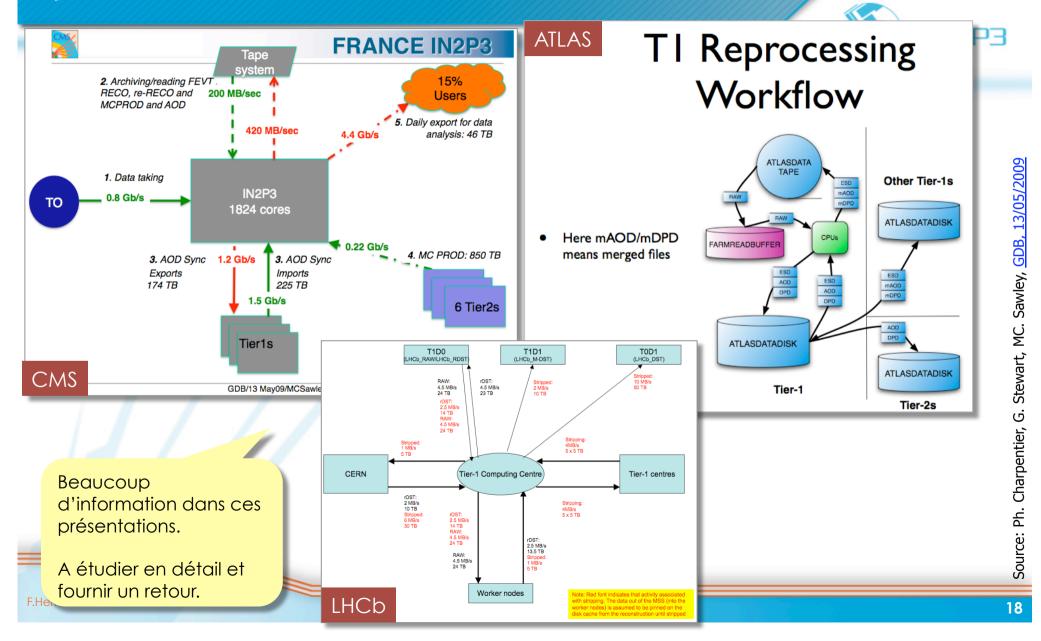
USER PID %CPU %MEM VSZ RSS TTY STAT START TIME COMMAND gla012 1741 0.1 0.0 58892 2612 pts/13 S 12:49 0:00 python yawn.py gla012 1745 0.0 0.0 5816 2176 pts/13 S 12:49 0:00 python32 yawn.py

 This means that running 'hello world' on the grid on a 64bit machine has a VMEM footprint of ~450MB (c.f. 32bit footprint of ~60MB)

Python 32bits VSZ = 5 MB

- In both cases the RSS is ~22MB
- Which means that killing grid jobs based on VMEM consumption is probably not a good idea...
 - N.B. also killing jobs based on RSS doesn't work as the kernel actively tries to keep pages in memory
 - And it seems that, e.g., torque kills based on the memory consumption of the process tree, not of the payload

Workflows & débits I/O





Tests glexec/SCAS



Status of the pilot

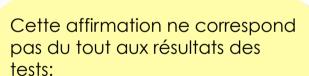
- Sites
 - FZK: Running well
 - Lancaster:
 - SCAS plug-in for 64 bit not building (ETICS issue). Will use pool account plug-in for now. Should be ready this week.
 - IN2P3: running well



- LHCb have started testing and are planning to increase this.
- ATLAS: Work on the framework is almost done. Expect to be sending jobs to prod in a few days.
- Issues:
 - Currently there is an incompatibility with CREAM (fix almost ready)



WLCG Grid Deployment Board, CERN 13 May 2009



- Impossibilité d'installer glexec partagé
- glexec détruit l'environnement du job lancé par le pilote

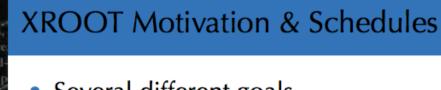
Détails Pierre Girard





XROOT comme protocol d'accès aux données





- Several different goals
- XROOT I/O server for CASTOR
 - Motivation:
 - more scalable I/O server for analysis
 - using XROOT server framework to implement a de-coupled CASTOR disk pool management (w/o LSF scheduling)
 - Time-scale: now (2.1.8)
 - aim at LSF free disk pools with next CASTOR release 2.1.9
- XROOT as potential common protocol
 - Motivation:
 - fewer file access protocols to support by users / providers
 - no intrinsic performance enhancements expected
 - Time-scale: unlikely that consolidation could finish before 2009/10 run
 - but planning / evaluation should start



CERNI"

Department

CERN IT Department CH-1211 Genève 23 Switzerland Dirk.Duellmann@cern.ch

Tuesday, 12 May 2009



Chantiers récents et en cours



- Migration vers HPSS v6.2
 - Interruption de service: Juin 1-4
- Pas de migration vers dCache/Chimera avant l'été
- Plusieurs tests et évolutions de la chaîne de soumission et de gestion des jobs
 - Voir présentation de Pierre aujourd'hui
- Problème d'installation et de réplication du logiciel ATLAS sur AFS
 - Réunion programmée le 29/05/2009



Chantiers récents et en cours (suite)



- Ménage dans les espaces de stockage (dCache et HPSS) des expériences LHC
- dCache
 - Poursuite de l'introduction de nouveaux serveurs de disques Thors et campagne de sortie de production progressive des premiers Thumpers (acquis en 2006)
- VO boxes
 - Toutes mise à jour et en production, à l'exception de une VO box ATLAS



RAPPEL: axes de travail prioritaires



- Tier-1
 - Améliorer l'efficacité de la copie bande → disque, via l'ordonnancement des requêtes passées à HPSS par dCache
- Séparation claire des données tier-1 et tier-2
 - Portes d'accès différentes pour chacun des sites, sans pour autant dupliquer les données
- Déploiement d'un prototype de ferme d'analyse interactif basée sur PROOF
- Finaliser le mécanisme d'installation du logiciel des expériences dans la zone AFS
 - Utilisation de la réplication de volumes AFS pour l'équilibrage de charge
 - Documentation et généralisation pour les 4 expériences LHC
- Monitoring !!!



Evénnements à venir



- Réunion semestrielle des sites LCG-France
 - Annecy, 18-19 mai 2009
 - http://indico.in2p3.fr/conferenceDisplay.py?confld=1660
- Hepix Spring 2009
 - Umeå (Suède), 25-29 mai 2009
 - http://www.hpc2n.umu.se/events/workshops/09/hepix/
- Workshop on adapting applications and computing services to multi-core and virtualization
 - CERN, Juin 24-26
 - http://indico.cern.ch/conferenceDisplay.py?confld=56353



Aujourd'hui et à venir



- Aujourd'hui
 - Evolutions chaîne de gestion de jobs
 - Etat d'avancement de l'interface CREAM pour BQS
 - SYMOD: plate-forme générique de monitoring
- Prochaine réunion
 - Jeudi 11 juin, 13h30-15h30, salle 202



Prochaines réunions en 2009



- Juin 11
- Juillet 9
- Septembre 10
- Octobre 15
- Novembre 12
- Décembre 10
- Toutes les réunions
 - http://indico.in2p3.fr/categoryDisplay.py?categld=102



Questions/Commentaires





F.Hernandez