



## Réunion des sites LCG-France

27-28 November 2008

LPSC, Grenoble

# ACTIVITÉS DU NUAGE ATLAS FRANÇAIS



**Excitement in the ATLAS Detector Control Room:  
The first LHC event on 10<sup>th</sup> September 2008**

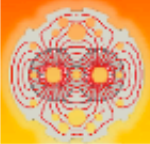
LHCC, 24-Sep-2008, PJ

2008...  
De Ahhh!!!!!!!!!!!!!!



**... as well as in the ATLAS Tier-0 and Data Quality Control Rooms:  
Reconstruction follow-up and analysis of the first LHC events**

LHCC, 24-Sep-2008, PJ



## Incident on 19th September



- During commissioning of the last main bend circuit to 5 TeV an incident occurred resulting in the triggering of quench heaters of about 100 magnets and a large He discharge into the tunnel.
- The most probable cause is a faulty electrical connection between two magnets. The sector is being brought to room temperature for repair.
- The time needed for warmup, repair and cooldown precludes a restart before CERN's obligatory winter shutdown.
- The shutdown schedule is being modified to gain ~ 1 month of LHC operation in 2009.

Lyn Evans - EDM5

A Ohhhh!!!!!!!!!!!!!!

2

# Un an et demi déjà...

4

## Conclusions

- Beaucoup de progrès en 2006
  - Transferts T0 → T1 → T2
  - Production MC distribuée T2 → T1 → T1
- A faire en 2007 :
  - Reprocessing (HPPS → Disque)
  - Mise en place de LCG3D (constantes de calibrations, etc..) tout ou presque est a faire
  - Analyse sur les T2 (rfio sur dpm non disponible...)
  - Finalisation du modèle d'analyse
    - Group-Analysis ~ bien défini
    - User-Analysis plus flou...
  - Amélioration des efficacités, outils de diagnostique
  - Final Dress Rehearsal : du T0 jusqu'a l'analyse dans les T2
- Composante française active dans ATLAS, cependant il faudrait :
  - Plus d'implication encore
  - Des liens plus étroits entre les sites





# 2008

5

- Année de stabilisation
  - ▣ Consolidation des outils
  - ▣ Amélioration des diagnostics
- Passage en mode opératoire
  - ▣ Tests systématiques (Functionnal Tests)
  - ▣ Shifts
  - ▣ Tickets
- Montée en puissance
  - ▣ Limitations observées
  - ▣ Mais.. Tjrs en dessous de ce que la réalité sera



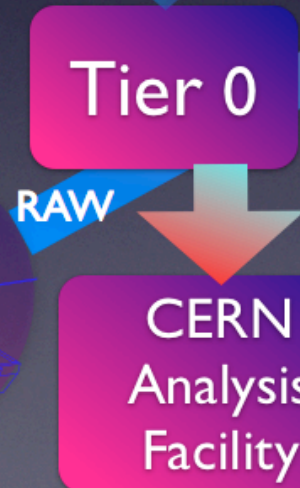
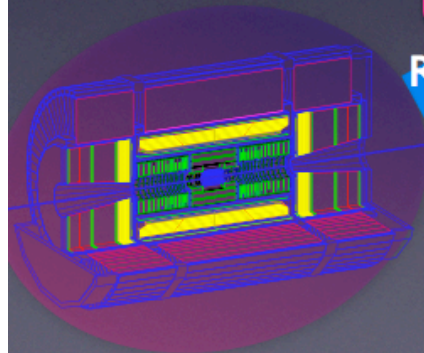
# The Computing Model

## Resources Spread Around the GRID

- Derive 1st pass calibrations within 24 hours.
- Reconstruct rest of the data keeping up with data taking.

- Reprocessing of full data with improved calibrations 2 months after data taking.
- Managed Tape Access: RAW, ESD
- Disk Access: AOD, fraction of ESD

- Interactive Analysis
- Plots, Fits, Toy MC, Studies, ...



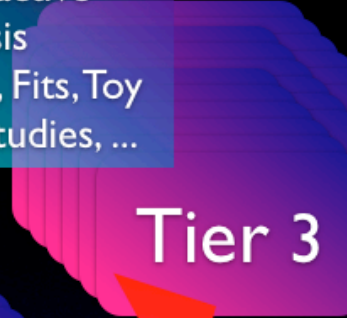
RAW/  
AOD/  
ESD



AOD



DPD



- Primary purpose: calibrations
- Small subset of collaboration will have access to full ESD.
- Limited Access to RAW Data.

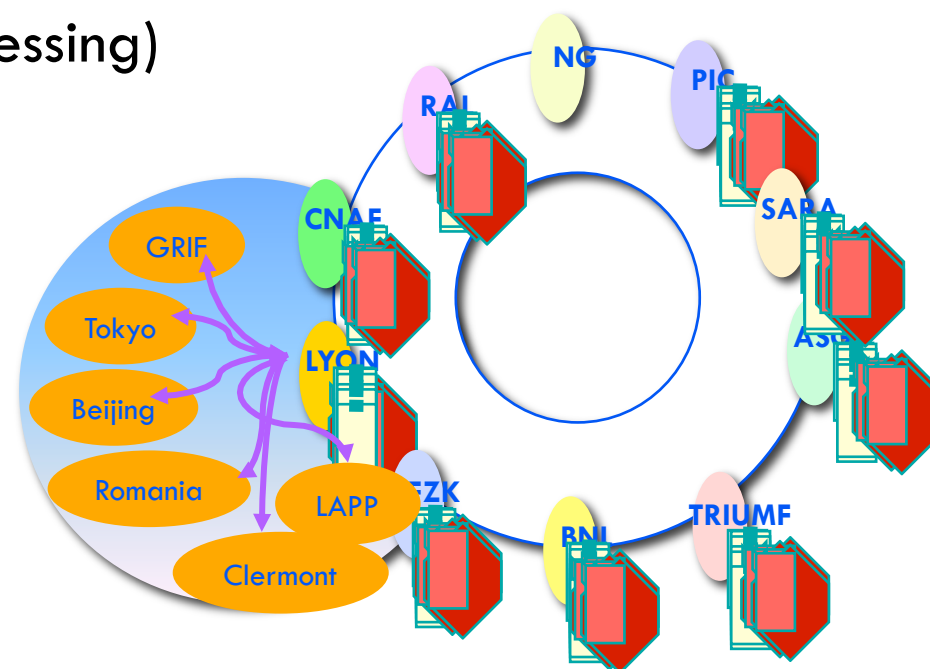
- Production of simulated events.
- User Analysis: 12 CPU/Analyzer
- Disk Store: AOD

© Amir Farbin/UTA  
CHEP 2007

# Le T1

7

- Stocke les données, en provenance
  - ▣ Du T0 (RAW & premier processing)
  - ▣ Des T1 (re-processing)
  - ▣ Des T2 (simulation MC)
- Distribue les données aux
  - ▣ T2s pour
    - Production MC (simulation)
    - Analyse
  - ▣ T1s pour archivage
- Process
  - ▣ Les données réelles (re-processing)
  - ▣ Les données MC (simulation & reconstruction)





# Le T1

Pierre angulaire du l'édifice (nuage)

8





# Les T2

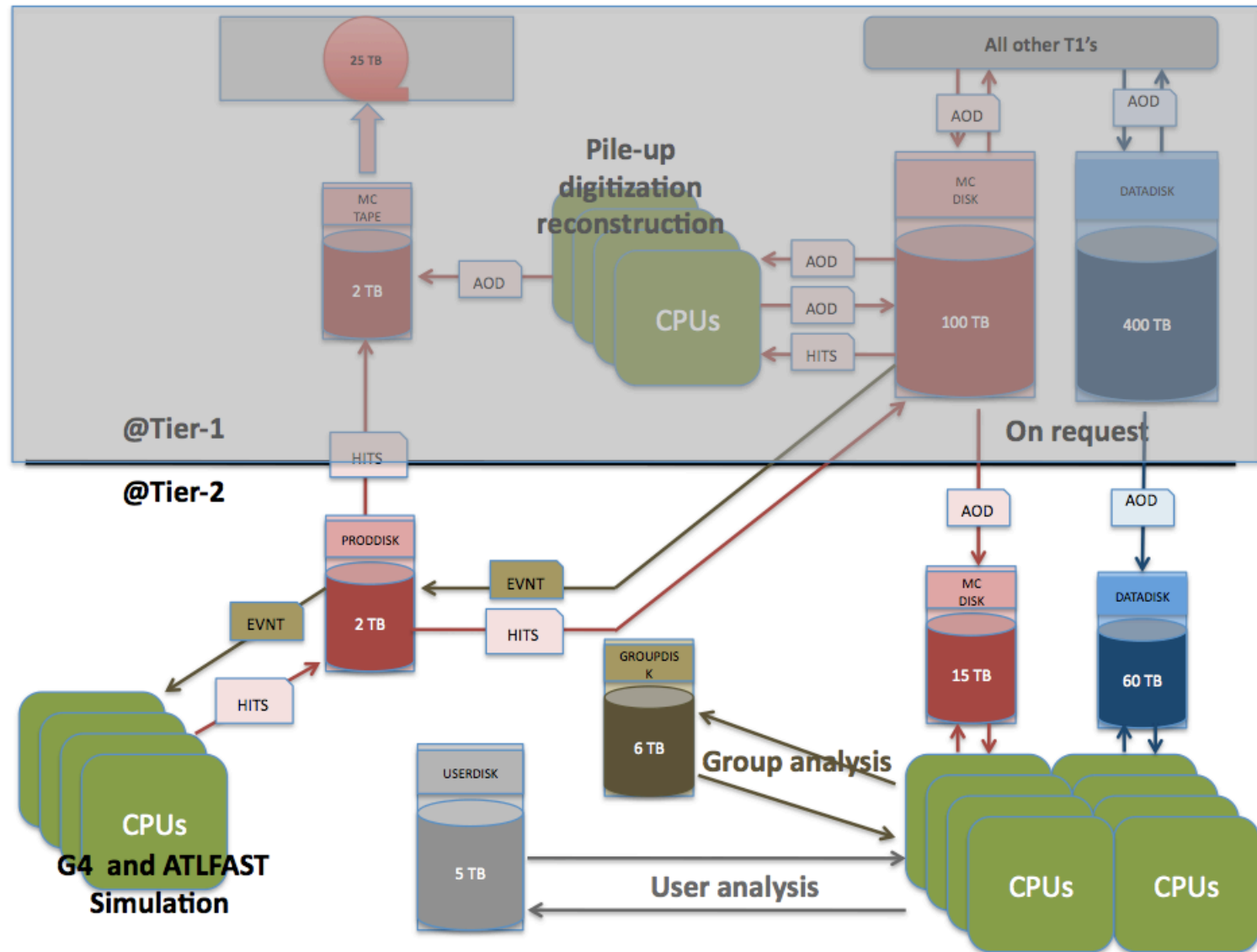
9

- Production MC
  - ▣ Envoyée au T1
- Analyse des données
  - ▣ Reçues du T1
- Pas d'échanges T2-T2 en dehors du nuage
  - ▣ Echanges inter nuages passent par le T1



# T1 & T2

10



# Problemes d'ATLAS

Sujets non abordes

11

- Evenements trop gros (x2)
- Temps de processing trop important (x1.5)
- Trop de PETITS fichiers
- Trop de memoire utilisee
- Modele de calibration non teste
- Modele d'analyse encore flou
- Duplication de certains outils
- Rationalisation des ressources
- ...

Ces problèmes ont  
des conséquences  
pour les sites



# Space Tokens

Nouveau en 2008

12

token name	storage type	used for	@T2	@T1	@T0
ATLASDATATAPE	T1D0	RAW data, ESD, AOD from re-proc		X	X
ATLASDATADISK	T0D1	ESD, AOD from data	X	X	X
ATLASMCTAPE	T1D0	HITS from G4, AOD from ATLFast		X	
ATLASMCDISK	T0D1	AOD from MC	X	X	X
ATLASPRODDISK	T0D1	buffer for in-and export	X		
ATLASGROUPDISK	T0D1	DPD de groupes	X	X	X
ATLASUSERDISK	T0D1	User Data	X	X *)	
ATLASLOCALGROUP DISK	T0D1	Local User Data @T3			

□ <https://twiki.cern.ch/twiki/bin/view/Atlas/StorageSetUp>

irfu

cea

saclay

Eric Lancon 27/11/08

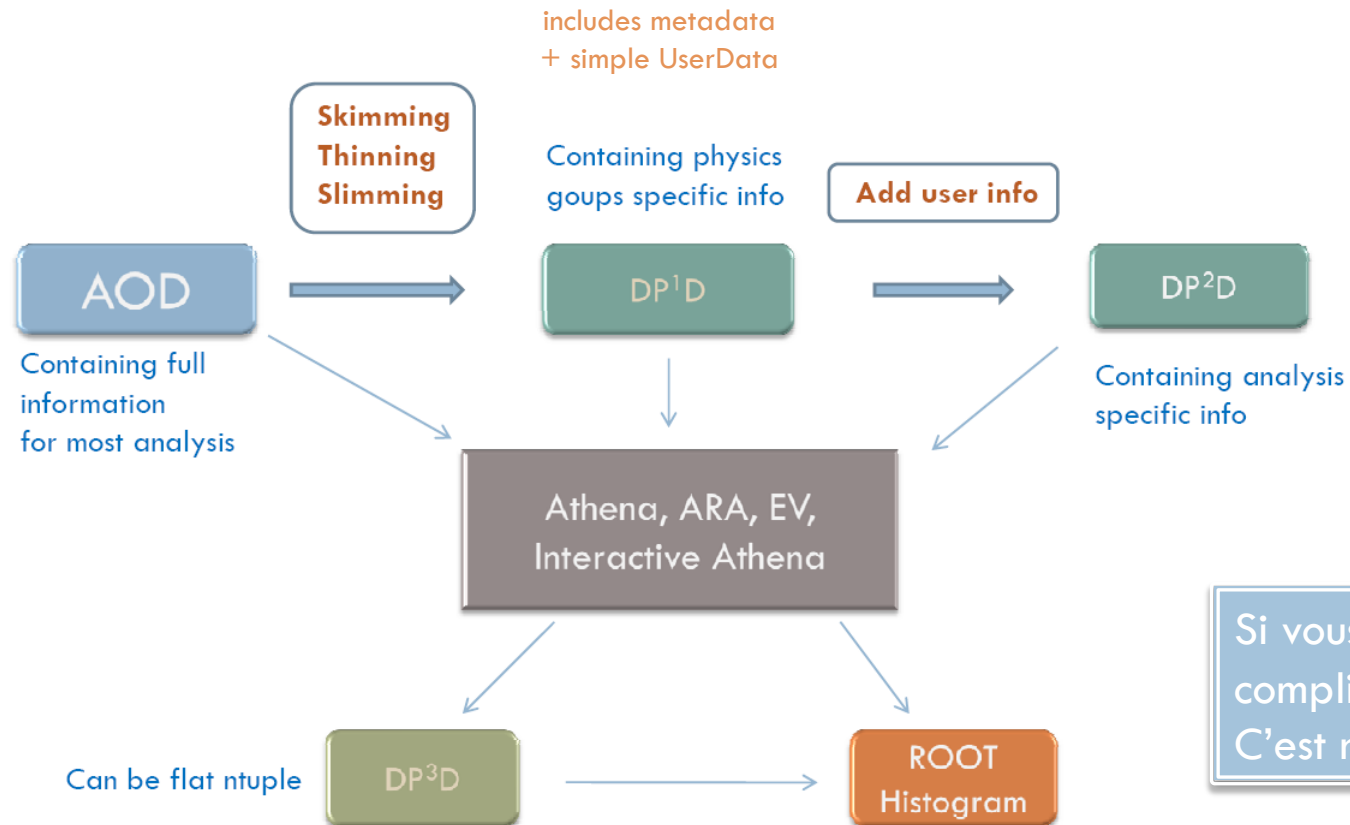




# Analysis Model

Nouveau en 2008

13



Si vous trouvez cela compliqué...  
C'est normal

DP1D centrally produced (T1), one per physic group (~10)



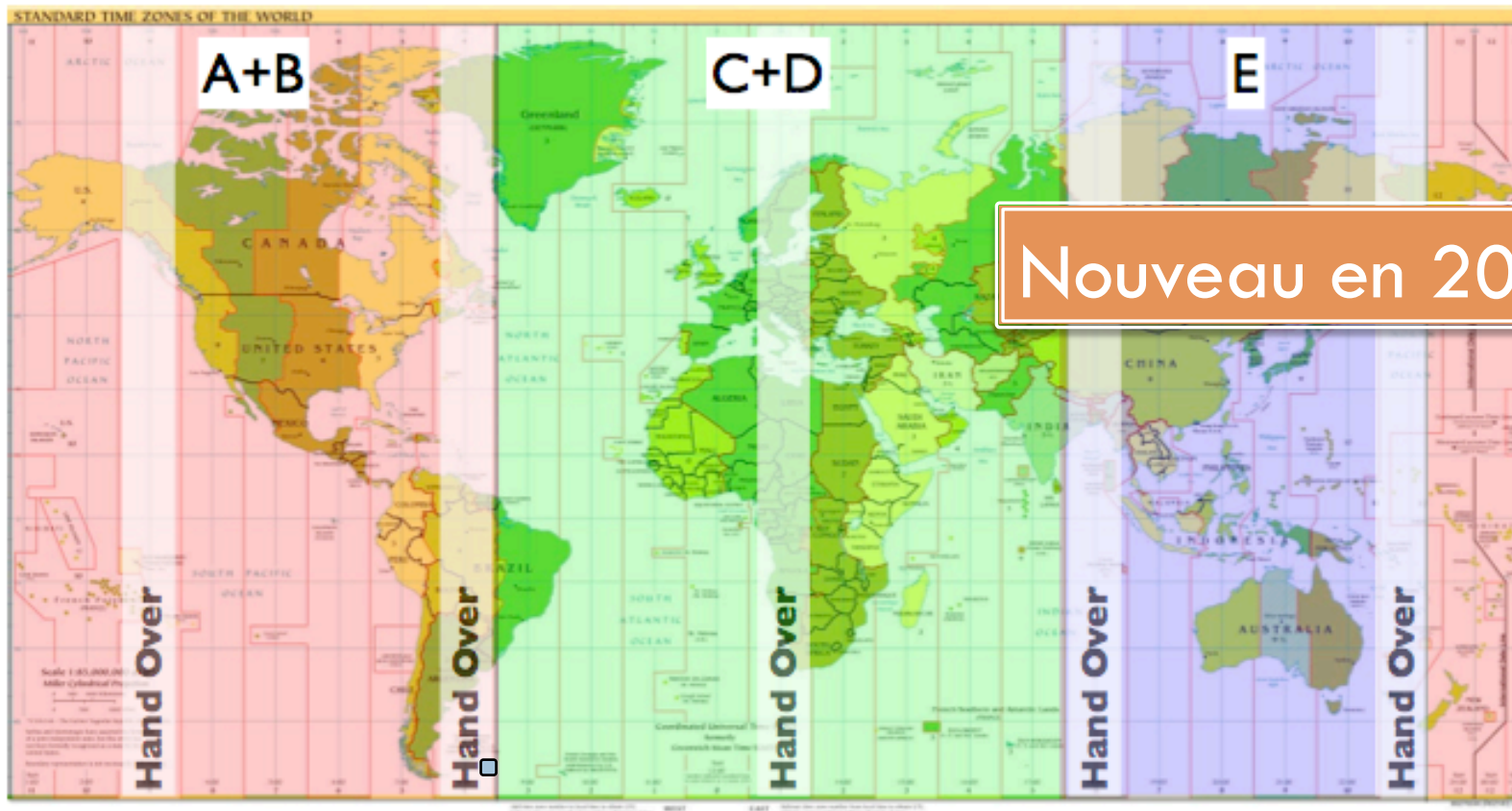
# ATLAS Distributed Computing in 2008



14

Full ADC operations and Round-the-clock shifts

Phase III: 5 Shifters on Duty (+ Trainees) - 24h coverage<sup>6x6h</sup>



irfu



saclay



# Stagein tests

Target : 1.2TB/h for a 10% T1  
Lyon problems  
dCache, HPSS

15

Nouveau en 2008

## Site Summary

Site	Rate GB/hr	Site	Rate GB/hr
CERN	2500	FZK	500
RAL	1000	LYON	700
TRIUMF	900	SARA	270
NDGF	200	CNAF	660
ASGC	2000	BNL	700
PIC	500(?)		

- Still many problems to resolve before we can have confidence in prestaging at T1s



# Data Distribution : Transfers to T1

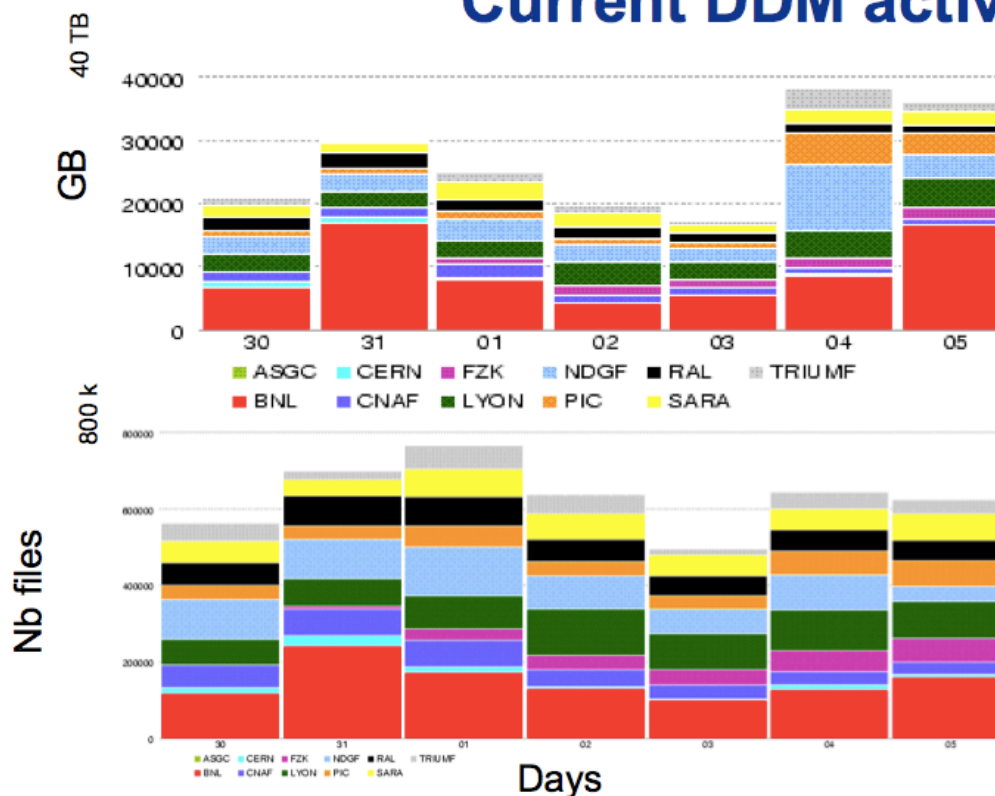
16

DDM : Data  
Distribution  
Management

Functional  
tests routinely  
performed

Some clouds  
clearly more  
efficient and  
more reactive  
than others

## Current DDM activity



Data placement/management

5

7 November 2009

irfu



saclay

Eric Lancon 27/11/08

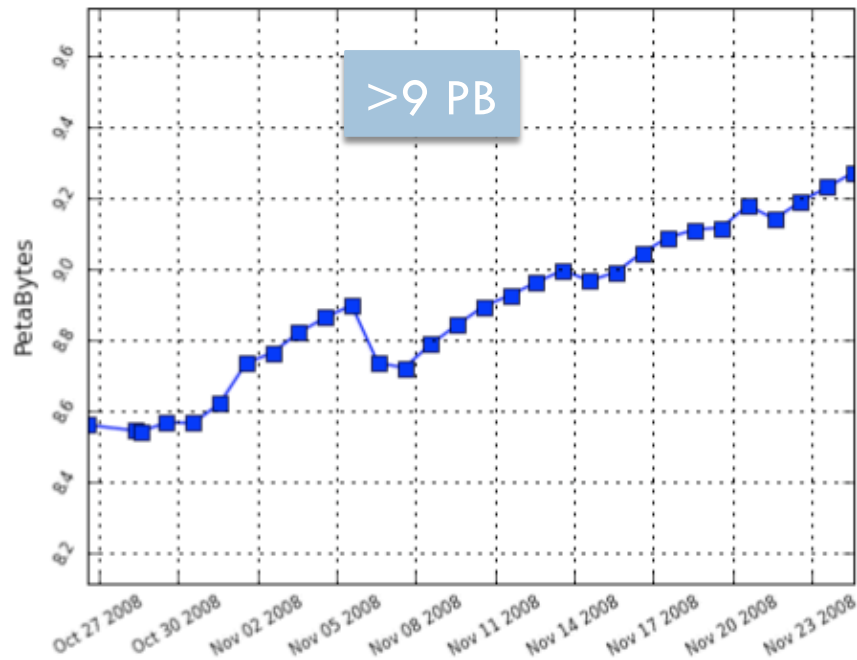




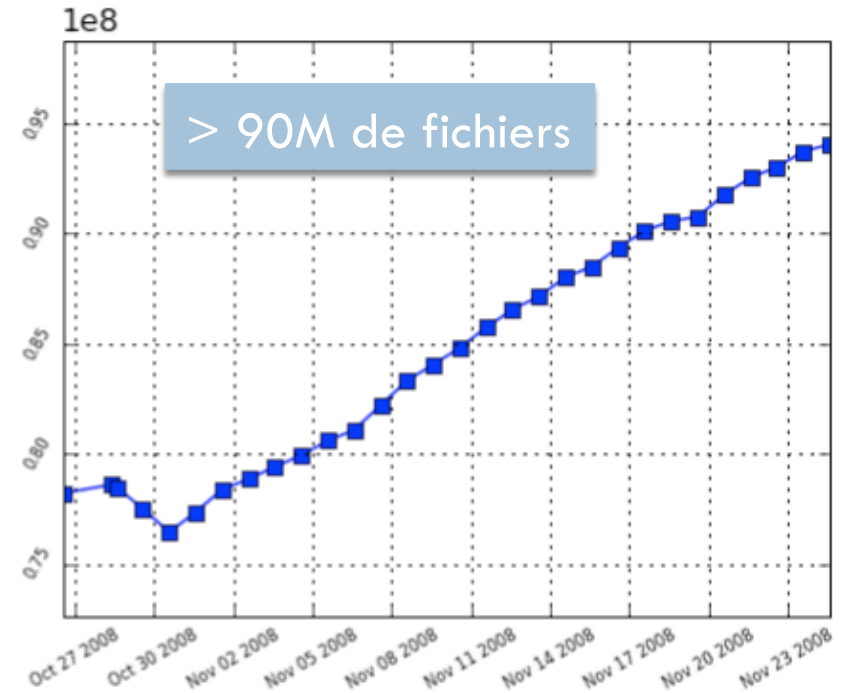
# Stockage ATLAS : Statistiques

17

Total GRID disk usage according to dq2



Total GRID files according to dq2



# Data replication in T2/3

18

□ [http://atladcops.cern.ch:8000/drmon/ftmon\\_TiersInfo.html](http://atladcops.cern.ch:8000/drmon/ftmon_TiersInfo.html)

□ For AODs & DPDs

□ If DPDs small enough

FRANCE	BEIJING-LCG2_DATADISK	100%	Egamma
	BEIJING-LCG2_DATADISK	35%	Minbias
	BEIJING-LCG2_DATADISK	50%	Muon
	GRIF-LAL_DATADISK	45%	
	GRIF-LPNHE_DATADISK	25%	
	GRIF-SACLAY_DATADISK	30%	
	IN2P3-CPPM_DATADISK	100%	Egamma
	IN2P3-LAPP_DATADISK	100%	Egamma
	IN2P3-LPC_DATADISK	50%	Egamma
	IN2P3-LPC_DATADISK	100%	Muon
	IN2P3-LPSC_DATADISK	5%	
	RO-02-NIPNE_DATADISK	10%	
	RO-07-NIPNE_DATADISK	10%	
	TOKYO-LCG2_DATADISK	100%	

□ Special replication request

□ Only for **Very special** cases...

FR-Cloud : 300% of AODs  
ATLAS model : 100%  
STRESS on T1 for distribution

irfu

cea

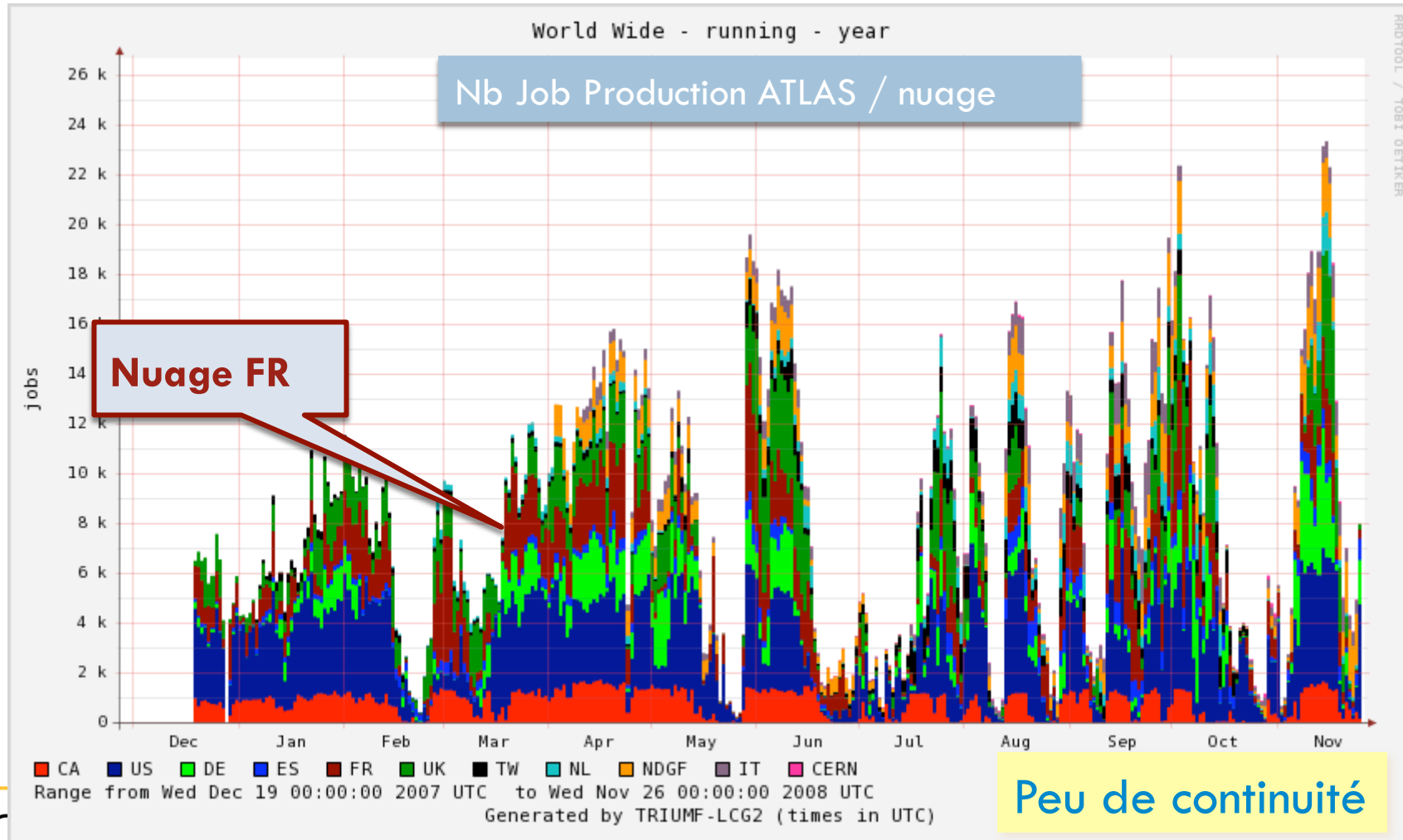
saclay

Eric Lancon 27/11/08



# Production MC

19



# Etat Production MC

20

## Conclusions

Trop de CPU

- More than enough CPU capacity for simulation needs Fall, 2008
- Major hurdle for recon production lowered (the vmem business)
- Throughput improved by more advanced task brokering
- Many of the clouds need to improve their efficiency and uptime dramatically and focus on storage
- We are progressing, but not clear if/when we will converge

Beaucoup de problèmes sur certains sites

irf

6 Nov, 2008

Production

17

cea

saclay

Eric Lançon 27/11/08





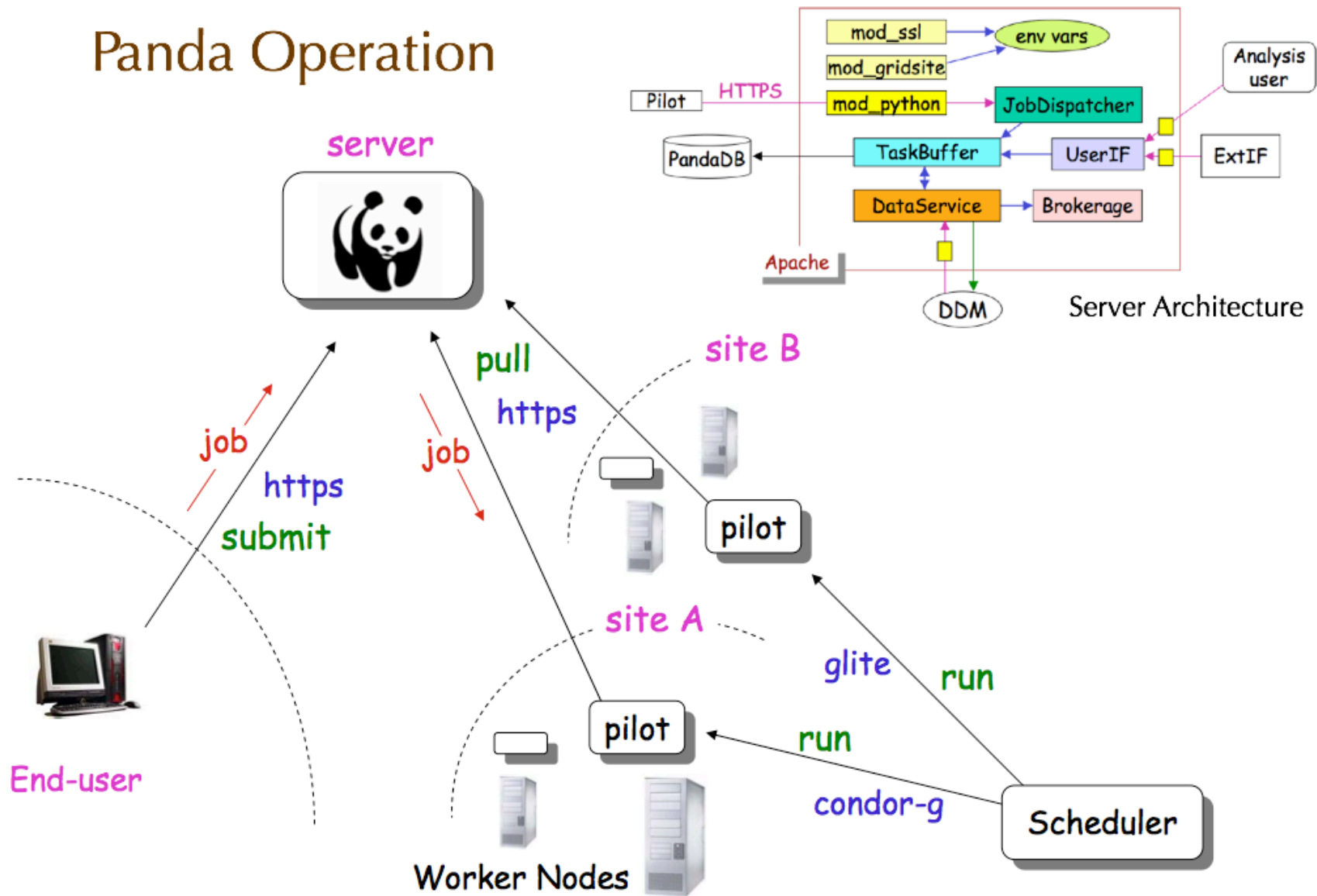
# Panda

21

- Outil unique utilise par ATLAS pour la production MC sur les trois types de grille
- A base de jobs pilots
  - ▣ Nous ne gâchons pas les ressources...

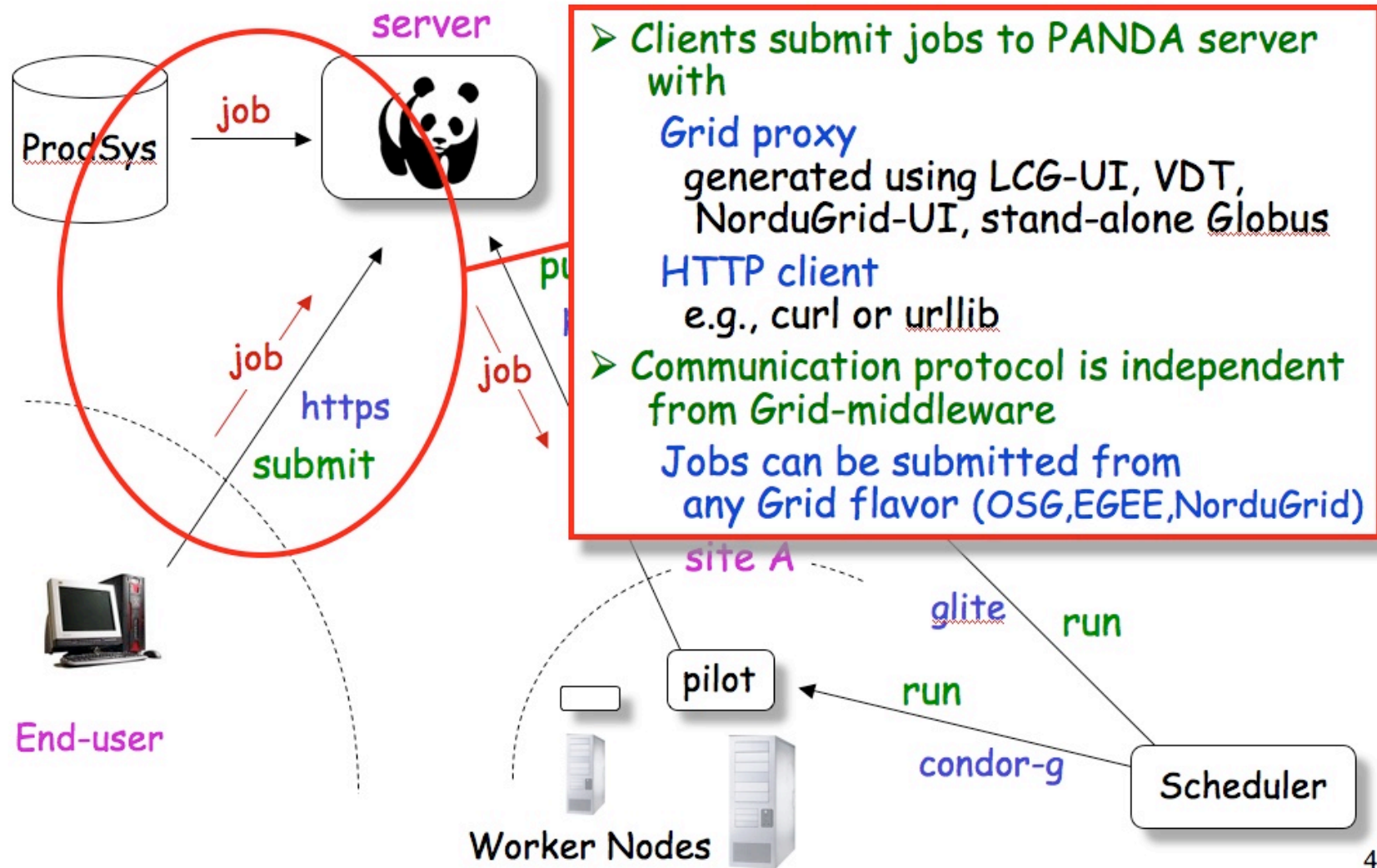


# Panda Operation





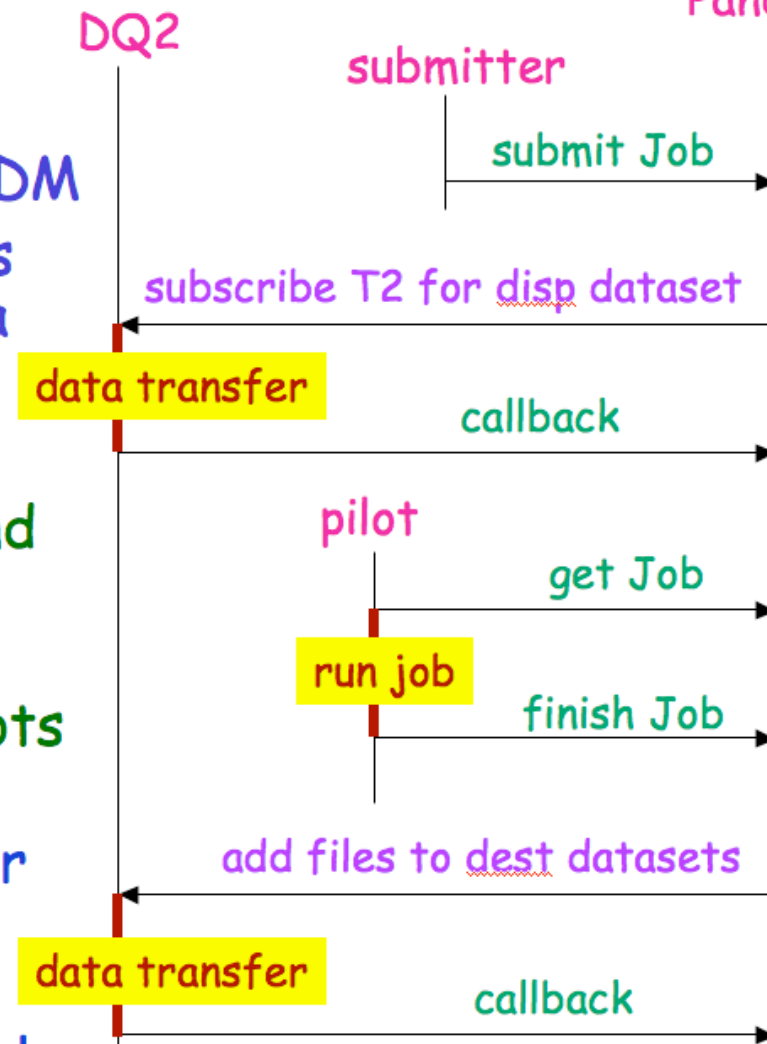
# Job Submission





# Data Transfer

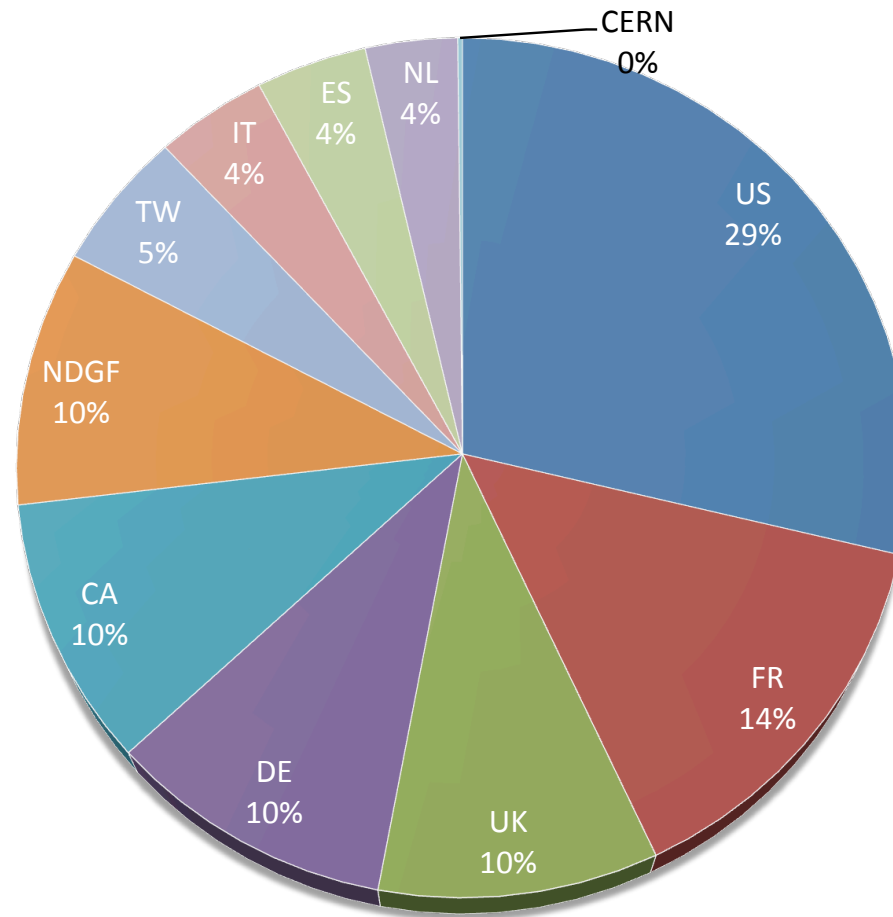
- Rely on ATLAS DDM
  - Panda sends requests to DDM
  - DDM moves files and sends notifications back to Panda
  - Panda and DDM work asynchronously
- Dispatch input files to T2s and aggregate output files to T1
- Jobs get 'activated' when all input files are copied, and pilots pick them up
  - Pilots don't have to wait for data arrival on WNs
  - Data-transfer and Job-execution can run in parallel



# Statistique par nuage

25

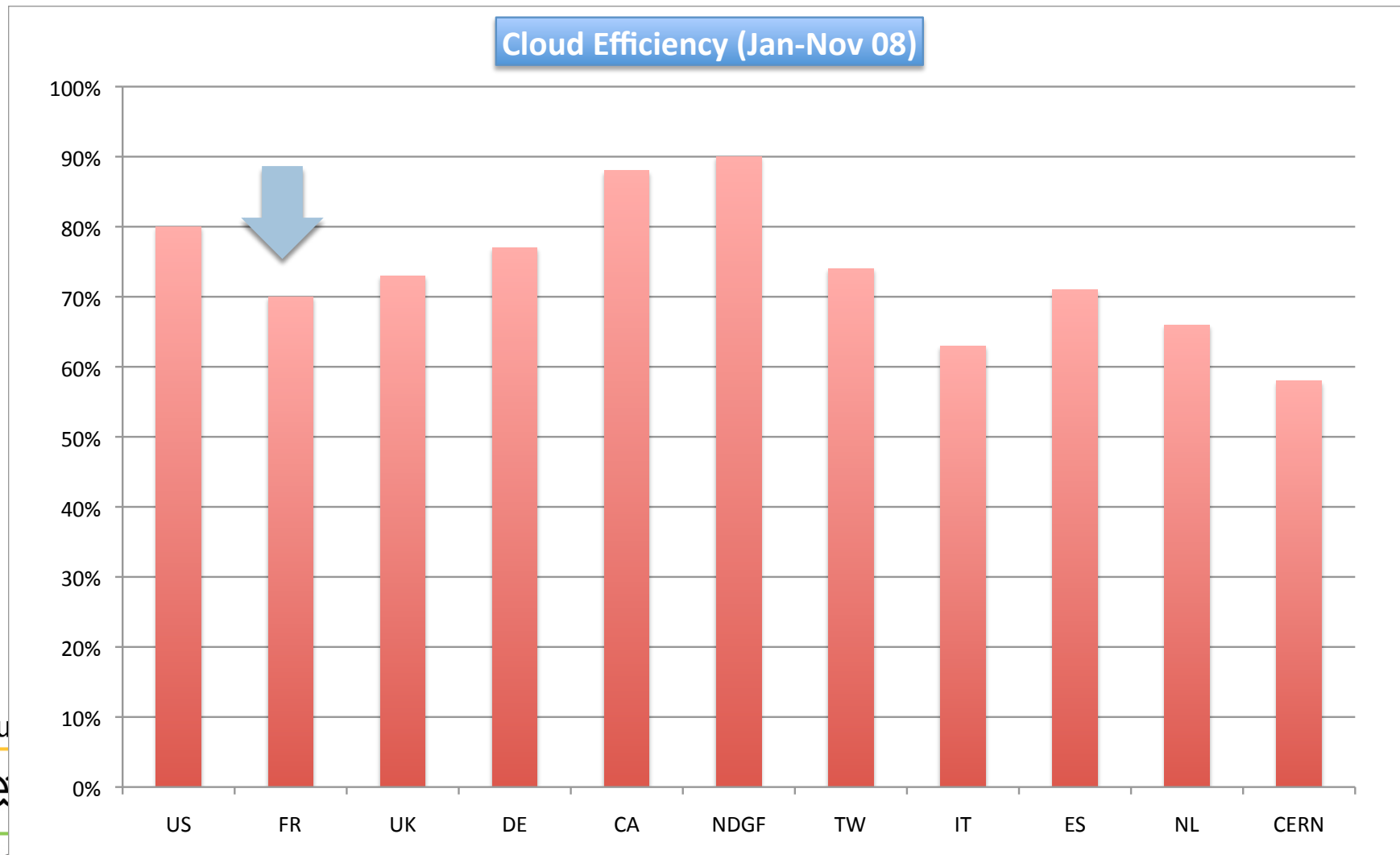
ATLAS MC production - Nb Jobs (Jan-Nov 08)





# Efficacité par nuage

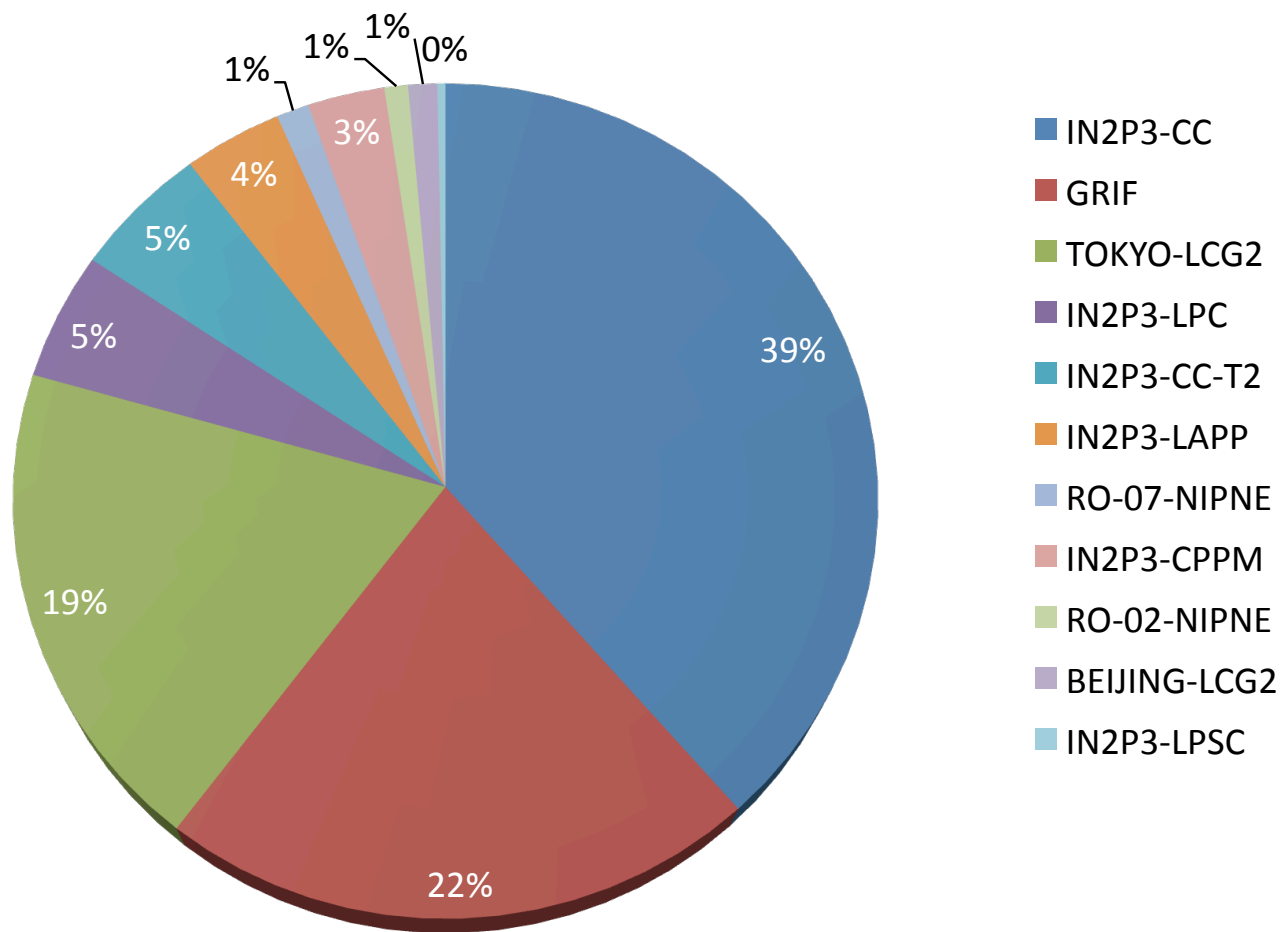
26



# Fraction de jobs par site du nuage

27

Fraction of jobs per site (Jan-Nov 08)



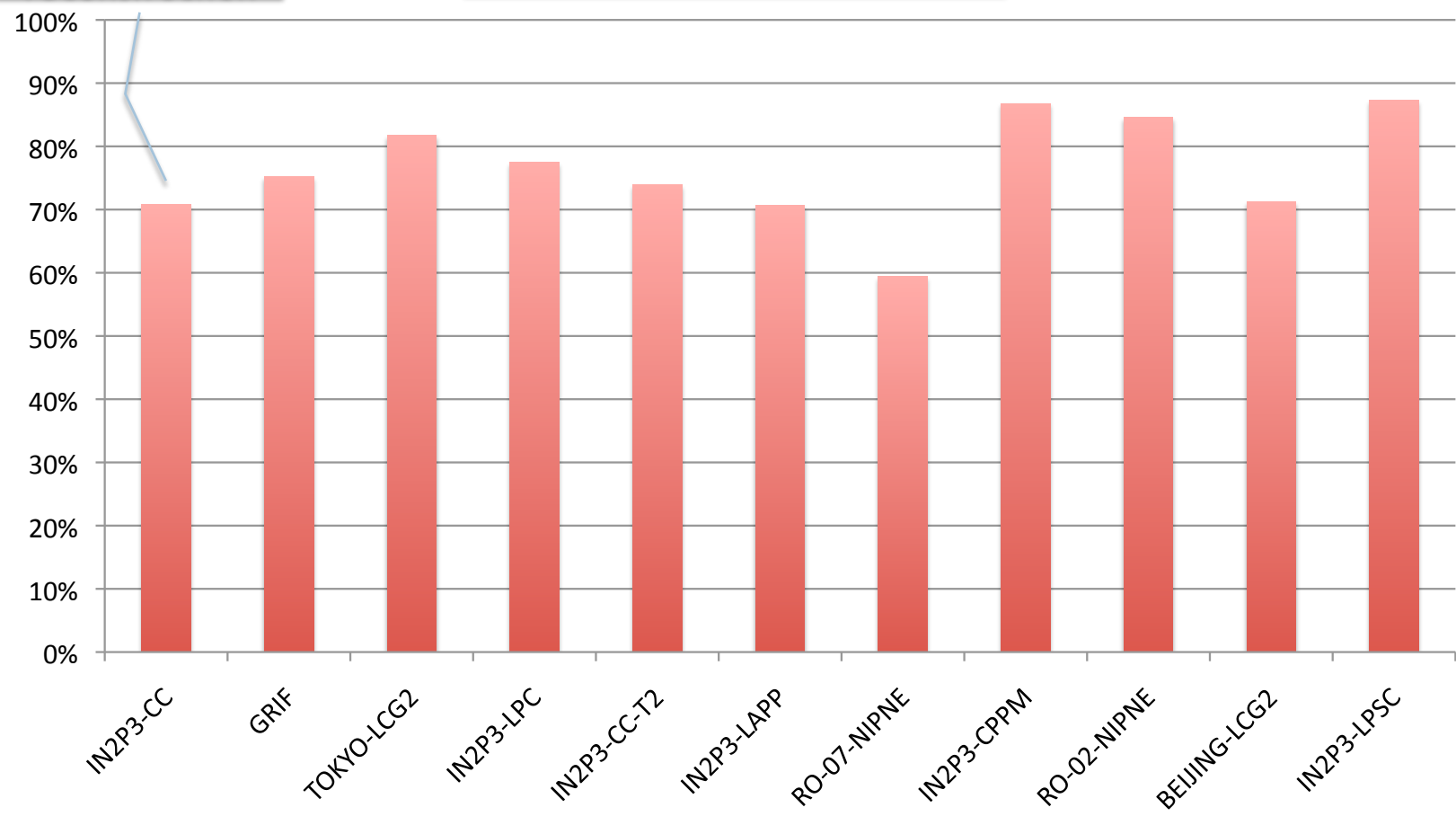
# Efficacité par site du nuage

28

Le seul à faire de la reconstruction

Par ordre d'importance en terme de jobs

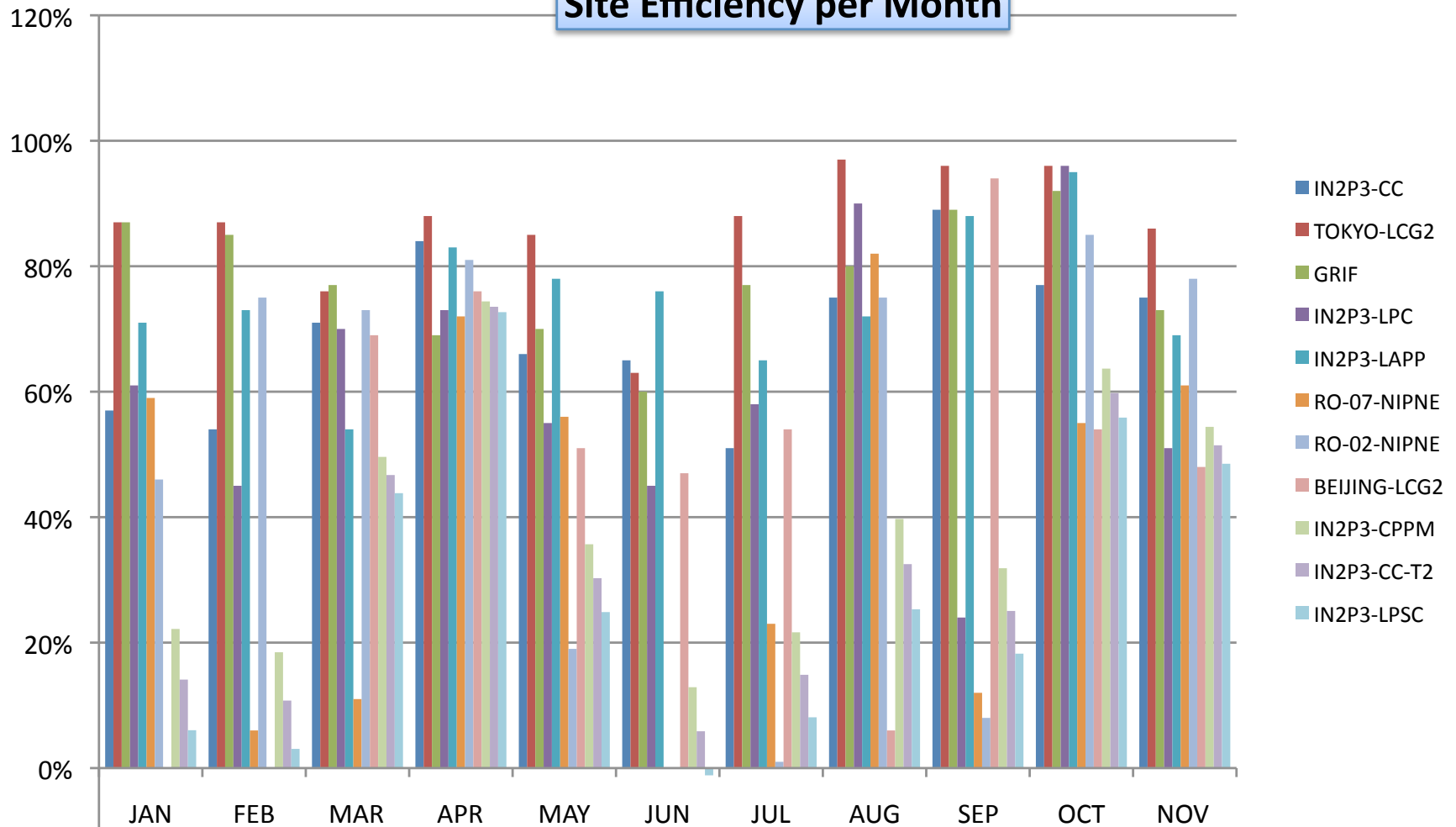
Job efficiency per site (Jan-Nov 08)



# Efficacité / mois

29

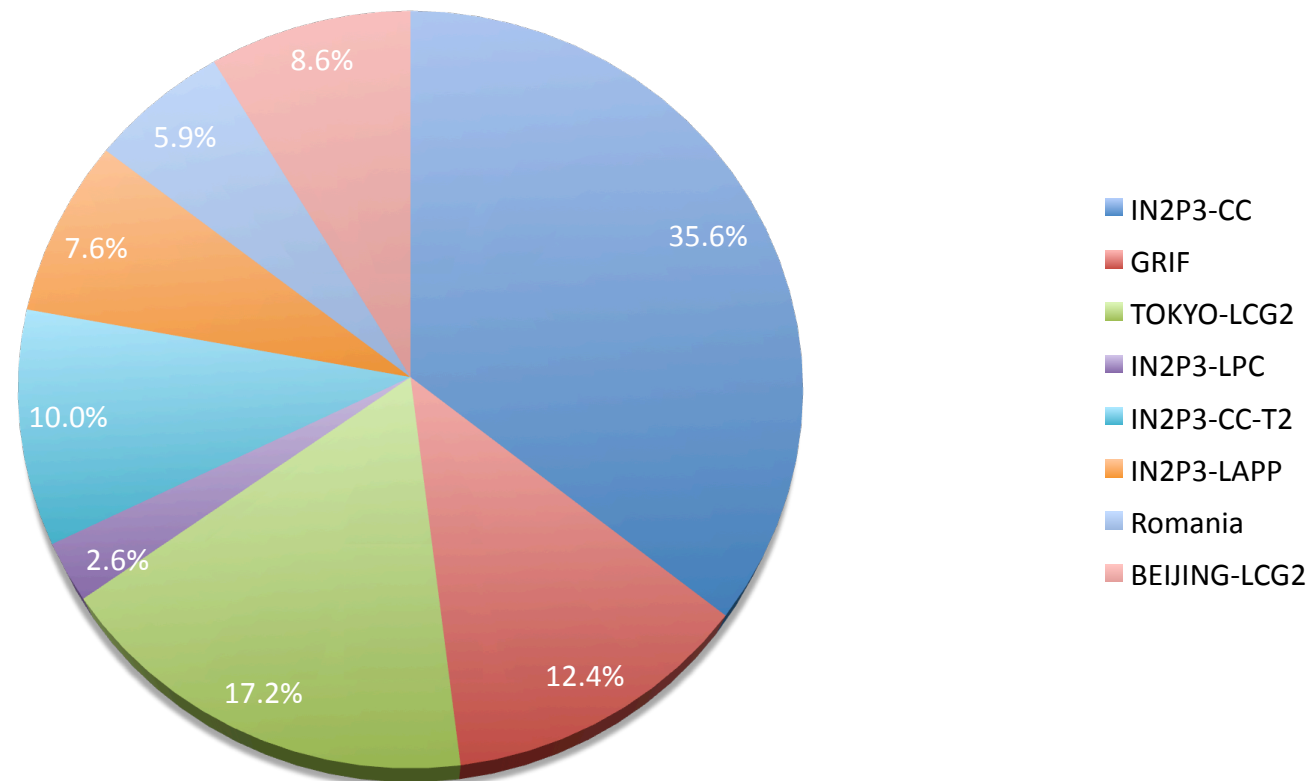
Site Efficiency per Month



# Pledges (CPU) per site

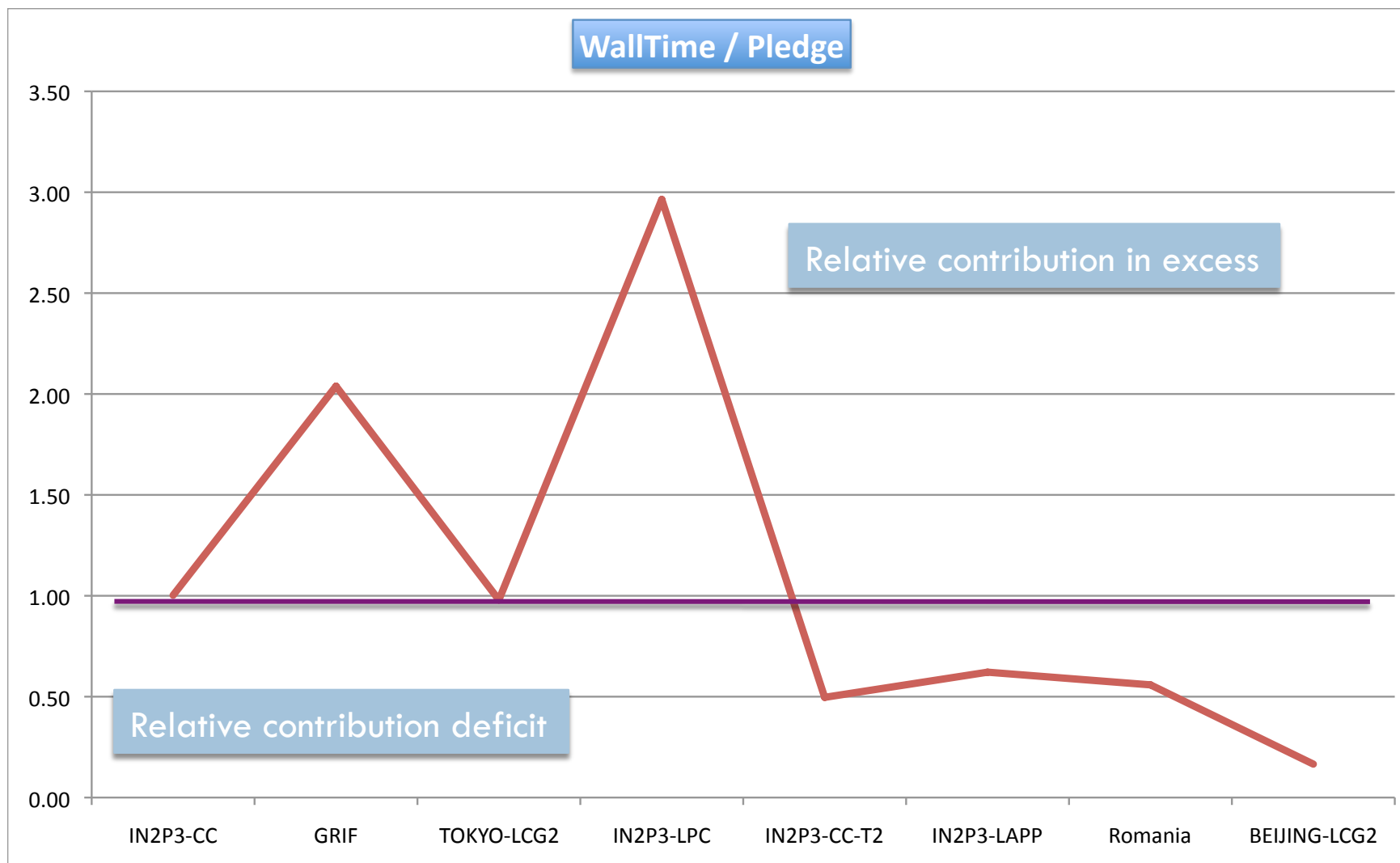
30

CPU pledges for 2008



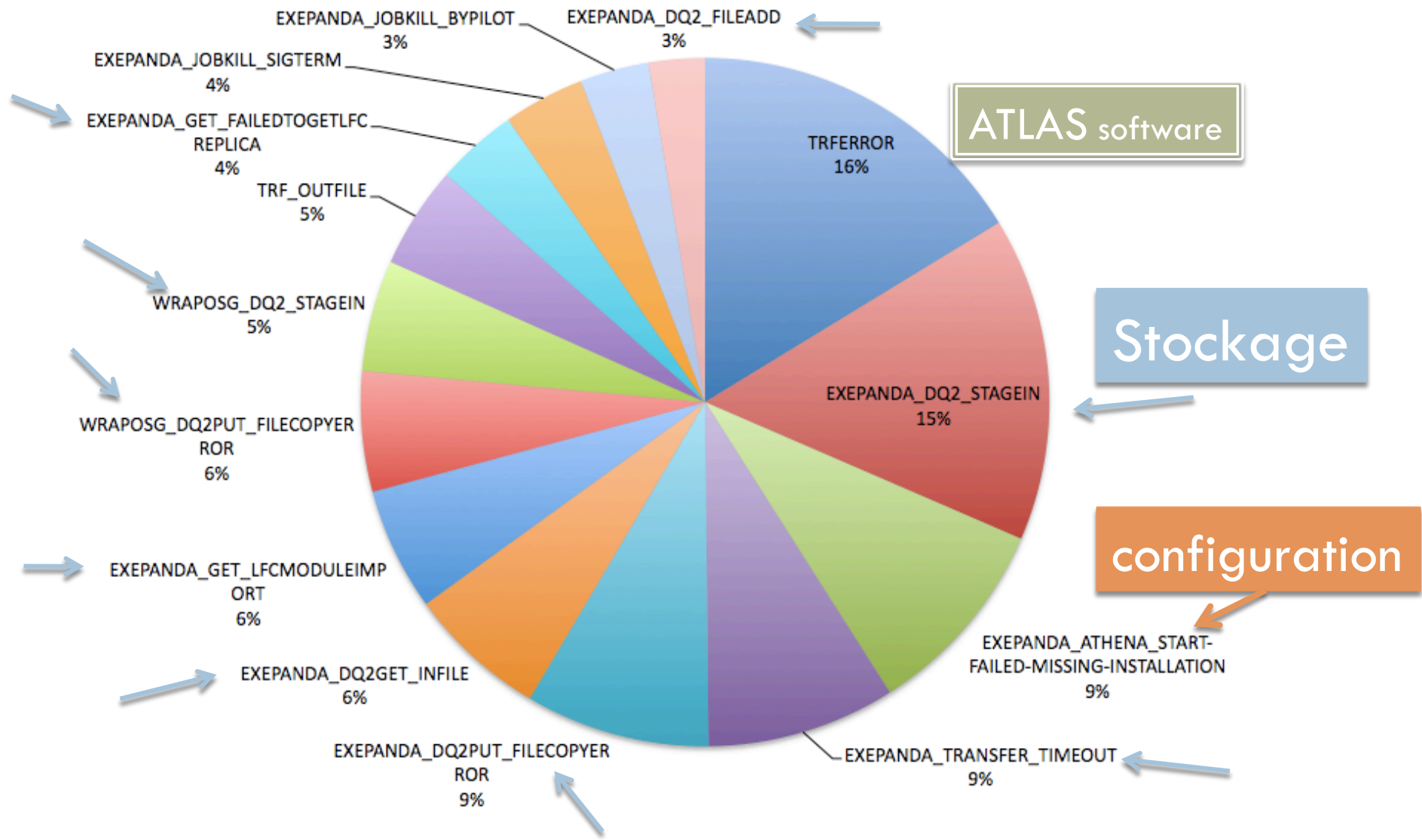
# Distribution relative des sites

31





# Most frequent errors (Jan-Nov 08)



# ATLAS au CC en 2008

33

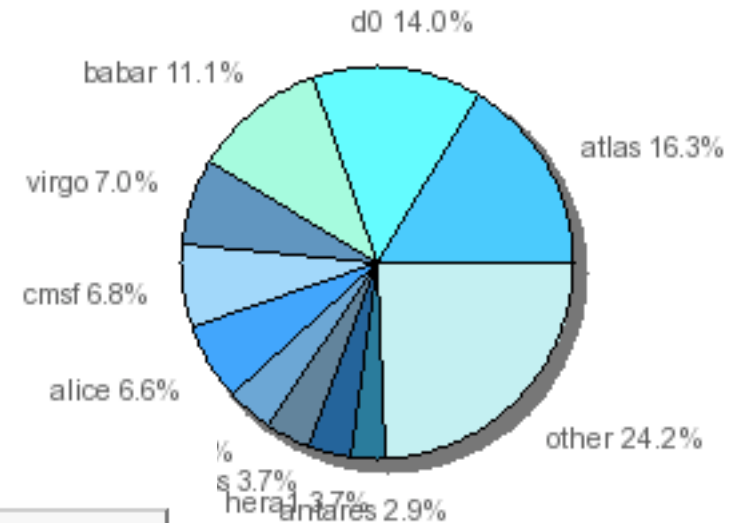
~25% des demandes utilisées

24/11/08

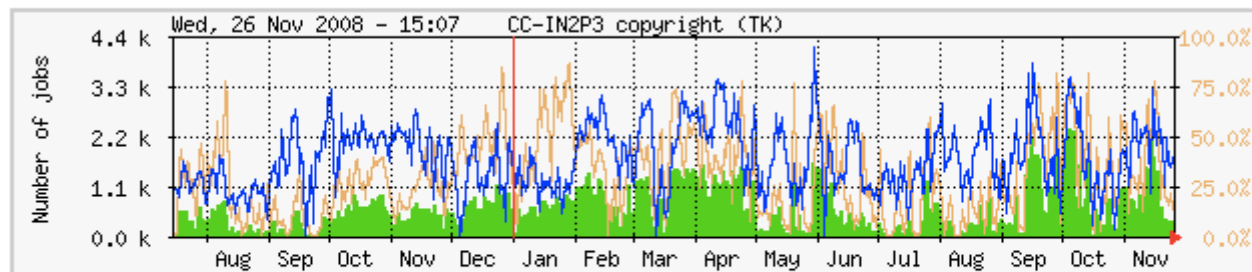
Details (order by consumption):

Range	Group	Topic	UI (hour)	Request for 2008	SI2k	%
1	atlas	lhc	163,124,395	565,000,000	1,036,105	16.28
2	d0	hep-nuc	140,358,969	230,000,000	891,508	14.01
3	babar	hep-nuc	111,674,001	185,000,000	709,311	11.14
4	virgo	astro-neutrino	70,414,720	40,000,000	447,248	7.03
5	cmsf	lhc	67,995,558	0	431,882	6.78

CC-IN2P3 Top 10 on anastasia farm



Yearly Graph (1 Day Average)



	Max	Average	Current
atlas grid running jobs:	2361	546	315
all grid running jobs:	4117	1693	1717
percentage:	86.0 %	32.0 %	18.0 %

saclay

Même en période de pic < objectifs  
Limitation imposée par le stockage



# Fast track cloud

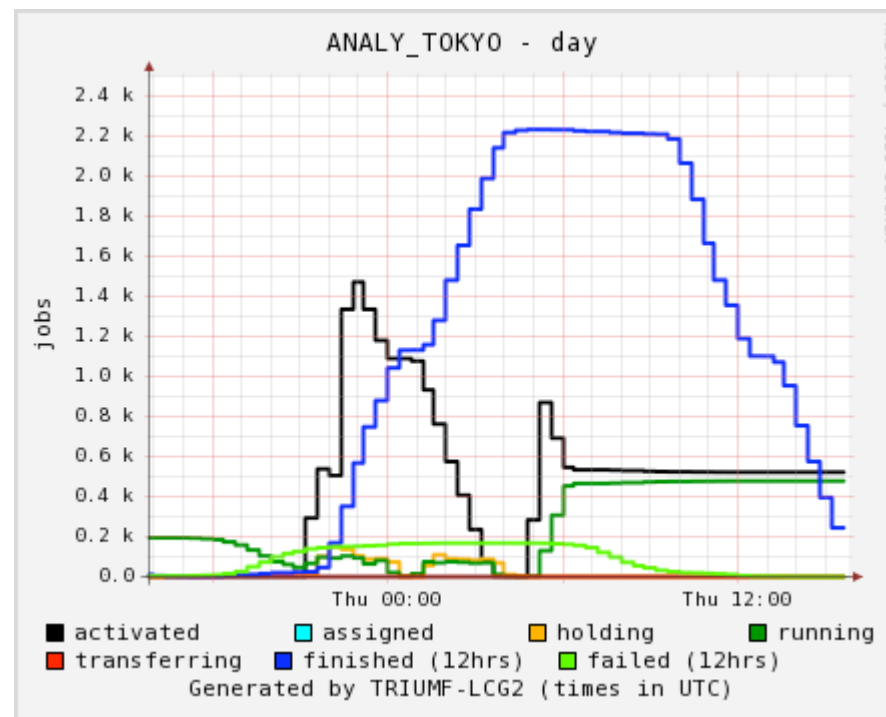
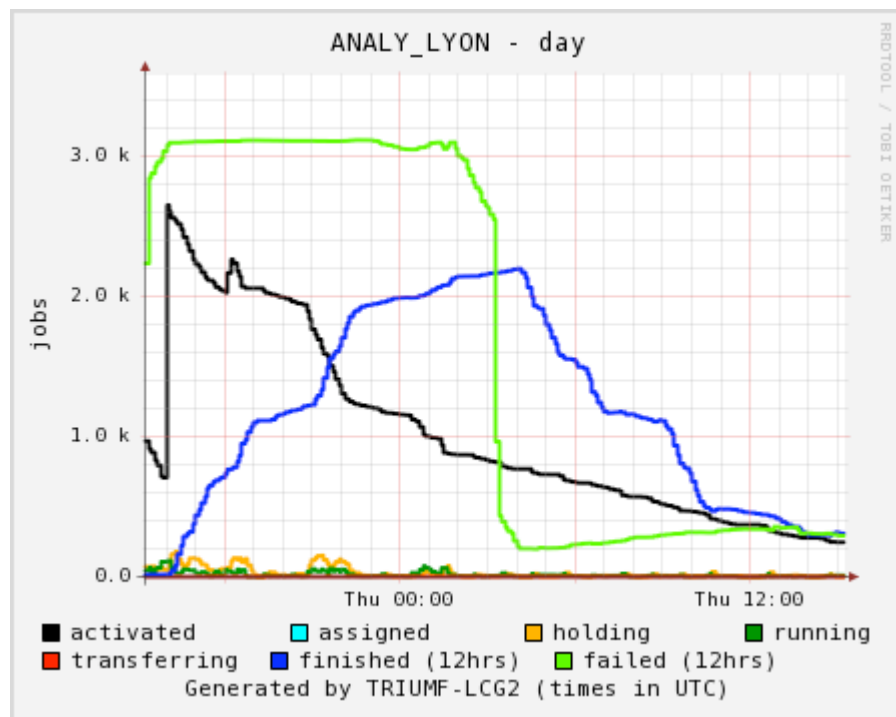
34

- Deux catégories de nuages pour ATLAS
  - ▣ 'ne pas s'embêter avec ceux qui posent problème'
  
- Fast Track Cloud (traduction : nuage fiable)
  - ▣ En priorité pour la production MC
  - ▣ Les tests de re-processing
  - ▣ ...
  - ▣ US, CA, NDF
  
- Malgré nos efforts pour améliorer l'efficacité du nuage FR nous sommes devenu un nuage de 'deuxième catégorie'
  - ▣ 2006 : FR un des 3 nuages de validations d'ATLAS
  - ▣ Trop de problèmes de stockage et distribution données at T1
  - ▣ Jobs ne s'exécutent pas assez vite



# Exemple de réactivité

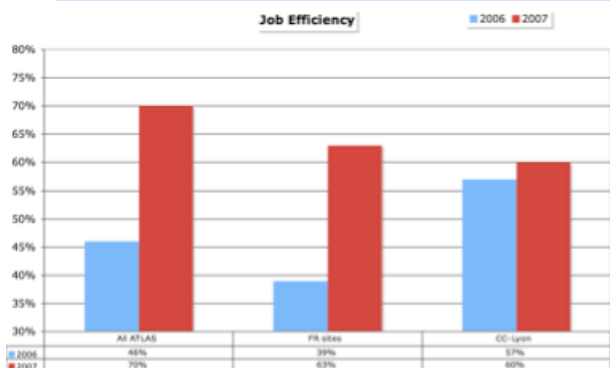
35



Exemple personnel : 3.000 jobs d'analyse dans (un peu plus que ) l'après midi au LPNHE , 3 jours au CC



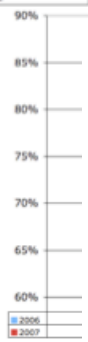
# .fr sites wrt others



- job efficiency
  - 2006 :
    - T1 & T2s better than <ATLAS>
  - 2007 :
    - T2s better than T1
    - but both worst than <ATLAS>

Warning présenté  
En début d'année  
Comite pilotage 30/1/08

- CPU efficiency
  - 2006 :
    - T1 & T2s better than <ATLAS>
  - 2007 :
    - T2s better than <ATLAS> ,
    - T1 worst than <ATLAS>



## 2007 MC production summary

- Big improvement of T2s (T3s) reliability
  - But Job efficiency lower than <ATLAS> because of storage instability at T1
- Degradation of T1 performances
  - Compared to both T2s
  - And <ATLAS>...
- Reasons of degradations :
  - /afs : affects only T1
  - dCache : affects all the cloud, MC production but data distribution as well



January-30-2008



## Current situation

- Deterioration of service quality at T1 over 2007
- It is not clear if with the current storage setup at Lyon, the ATLAS T1 could meet the following simultaneous requirements :
  - Storage of Data from CERN
  - Storage of MC Data from T2/T3
  - Data server for Analysis at T2 and CC-Lyon
- There is a real risk of missing the goals
  - No continuous Data flow from CERN tested
  - No sustained analysis yet (but french physicists tendency is to go to BNL for analysis...)

Les T2s sont/seront aussi affectes car les données sont distribuées depuis le T1

Les tests d'analyse (par exemple) sont faits sur avec des données pré positionnées

Ce qui n'est pas le mode opératoire en période de prise de données

Physicien  $\Lambda$  : comment le candidat H n'est pas arrive sur mon T2...

Comite pilotage 30/1/08



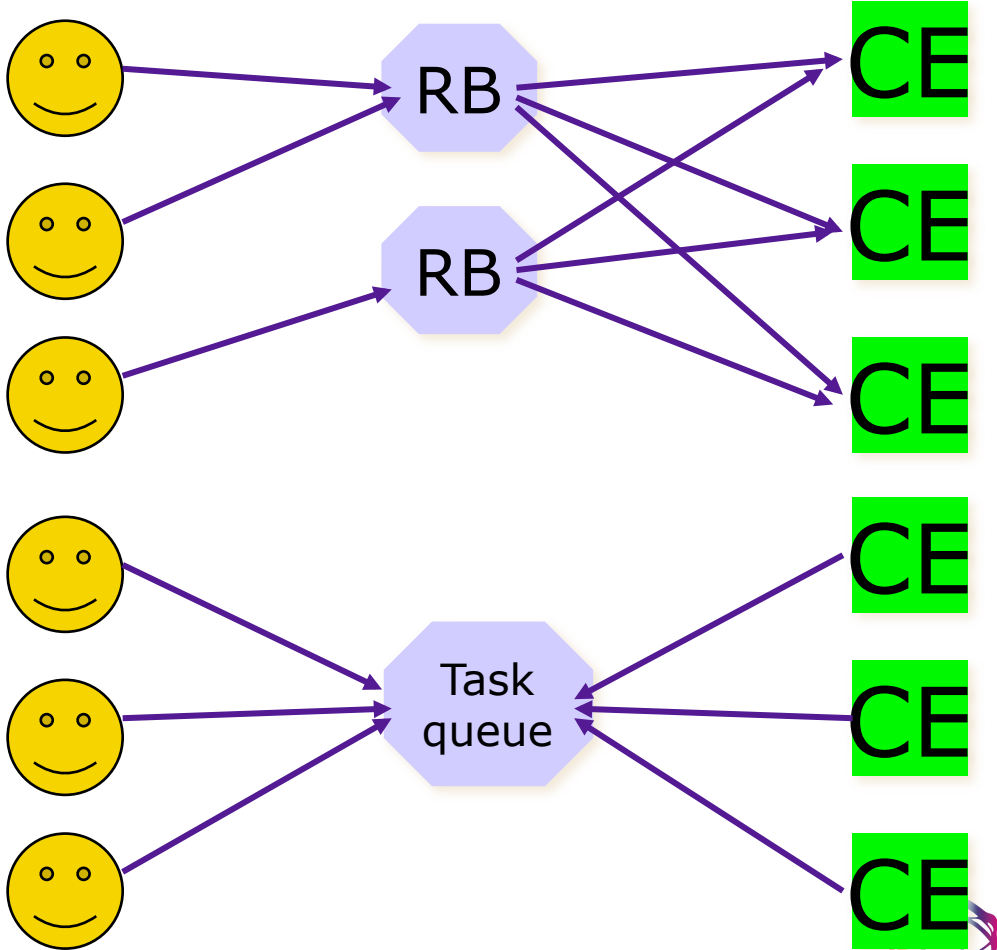
## Conclusions

- French T2 & T3 are now active players
- Tier-1 : Degradation of service
  - Too few people on key services (/afs, storage),
  - Very good will but overloaded and some problems of continuity during vacations etc...
  - Rapid actions have to be taken to
    - Investigate improvements/alternatives
    - Increase support
  - There is a real risk to not be ready.
  - Work is ahead of us... (see next page)



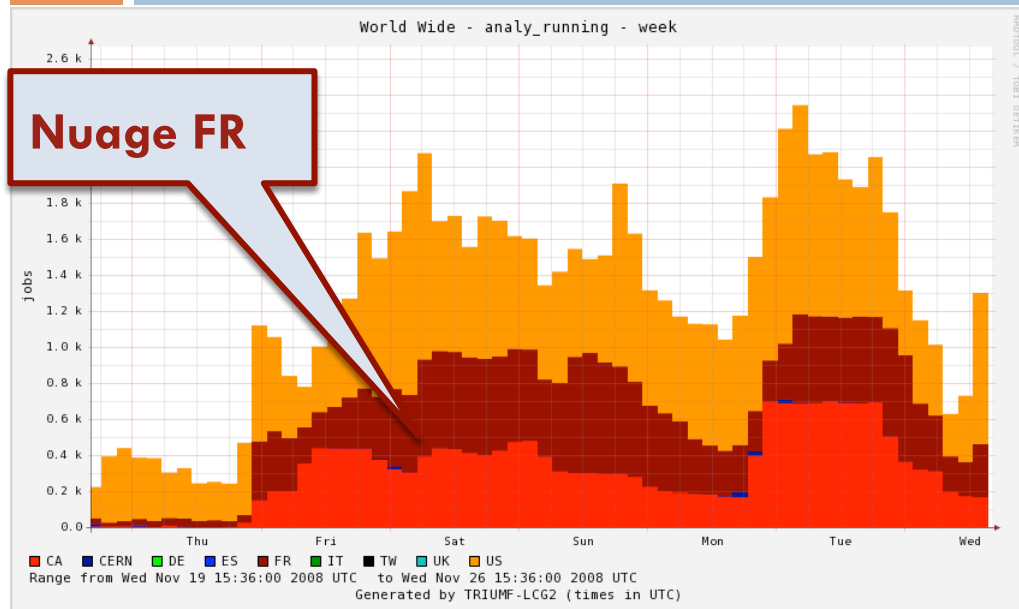
# Distributed analysis tools

- Ganga
  - <http://cern.ch/ganga>
  - Developed on LCG
  - Push mode
  - Only for analysis
  
- pAthena
  - Developed on OSG
  - Pull mode (Pilot jobs)
  - Fully integrated with DQ2
    - Outputs (libraries, logs, Ntup, etc..) are DQ2 datasets accessible through Panda browser
  - Same tool for
    - Production
    - Data analysis



# pAthena sur FR

39



Nuage FR

De plus en plus d'utilisateurs NON-FR

Viennent aussi des US!

Est ce que BNL restera toujours l'option 'facile' pour les français qui y vont?

Pas seulement a Lyon

st	Pilots (3hrs)	defined	assigned	waiting	activated	sent	running	holding	transferring	finished	failed	tot	trf	other
12:24	1723	201	0	0	422	0	222	65	0 / 0	709	467	40%	2%	38%
<a href="#">ANALY_BEIJING</a>	2	2	11-24 18:50	63	0	0	0	0	0 / 0	0	2	100%	0%	100%
<a href="#">ANALY_CPPM</a>	2	0	11-24 18:57	72	0	0	0	0	0 / 0	2	0	0%	0%	0%
<a href="#">ANALY_GRIF-IRFU</a>	2	0	11-24 18:43	87	0	0	0	0	0 / 0	2	0	0%	0%	0%
<a href="#">ANALY_GRIF-LAL</a>	2	0	11-24 18:57	94	0	0	0	0	0 / 0	2	0	0%	0%	0%
<a href="#">ANALY_GRIF-LPNHE</a>	16	0	11-24 22:23	243	0	0	199	2	25	0 / 0	78	0	0%	0%
<a href="#">ANALY_LAPP</a>	51	35	11-24 22:23	343	101	0	0	114	0	0 / 0	26	35	57%	0%
<a href="#">ANALY_LONG_LYON</a>	171	29	11-24 22:23	209	0	0	0	75	0	0 / 0	81	29	26%	0%
<a href="#">ANALY_LPC</a>	20	0	11-24 22:23	102	0	0	174	6	16	0 / 0	8	0	0%	0%
<a href="#">ANALY_LPSC</a>	0	0		109	0	0	0	0	0 / 0	0	0			
<a href="#">ANALY_LYON</a>	174	101	11-24 22:24	191	0	0	49	25	23	0 / 0	132	101	43%	9%
<a href="#">ANALY_LYON_XROOTD</a>	5	0	11-24 20:16	26	0	0	0	0	0 / 0	5	0	0%	0%	0%
<a href="#">ANALY_ROMANIA02</a>	0	0		32	0	0	0	0	0 / 0	0	0			
<a href="#">ANALY_ROMANIA07</a>	3	298	11-24 22:05	34	0	0	0	0	0 / 0	0	298	100%	0%	100%
<a href="#">ANALY_TOKYO</a>	117	2	11-24 22:19	118	100	0	0	0	1	0 / 0	373	2	1%	0%

irfu



saclay



# Analysis Challenge

40

- Performed on IT & DE clouds
  - ▣ Ganga only
  - ▣ No MC production
- FR-cloud next round
  - ▣ Ganga & pAthena
  - ▣ + MC production
- Interactions avec les sites
  - ▣ Pour Monitoring
  - ▣ Quels sites? (le plus possible!)
  - ▣ Frédérique Chollet (LGG-FR T2/3) contact pour les sites
- Soumissions de jobs
  - ▣ Coordonnée
  - ▣ Jobs réels!
- Quand ?  $\geq$  8 décembre



# IT – Cloud Analysis Challenge

- Started at 10:00 CET 27 October 2008
  - Used Ganga 5.0.10 + AthenaSplitterJob class modified to perform dataset-wise splitting.
  - Combination of individual & bulk WMS submissions
  - 700-800 jobs submitted to each site
- Results taken after 48 hours.

- 4 sites
- ~400 job/site/jour
  - ▣ C'est peu...



# Efficacités

42

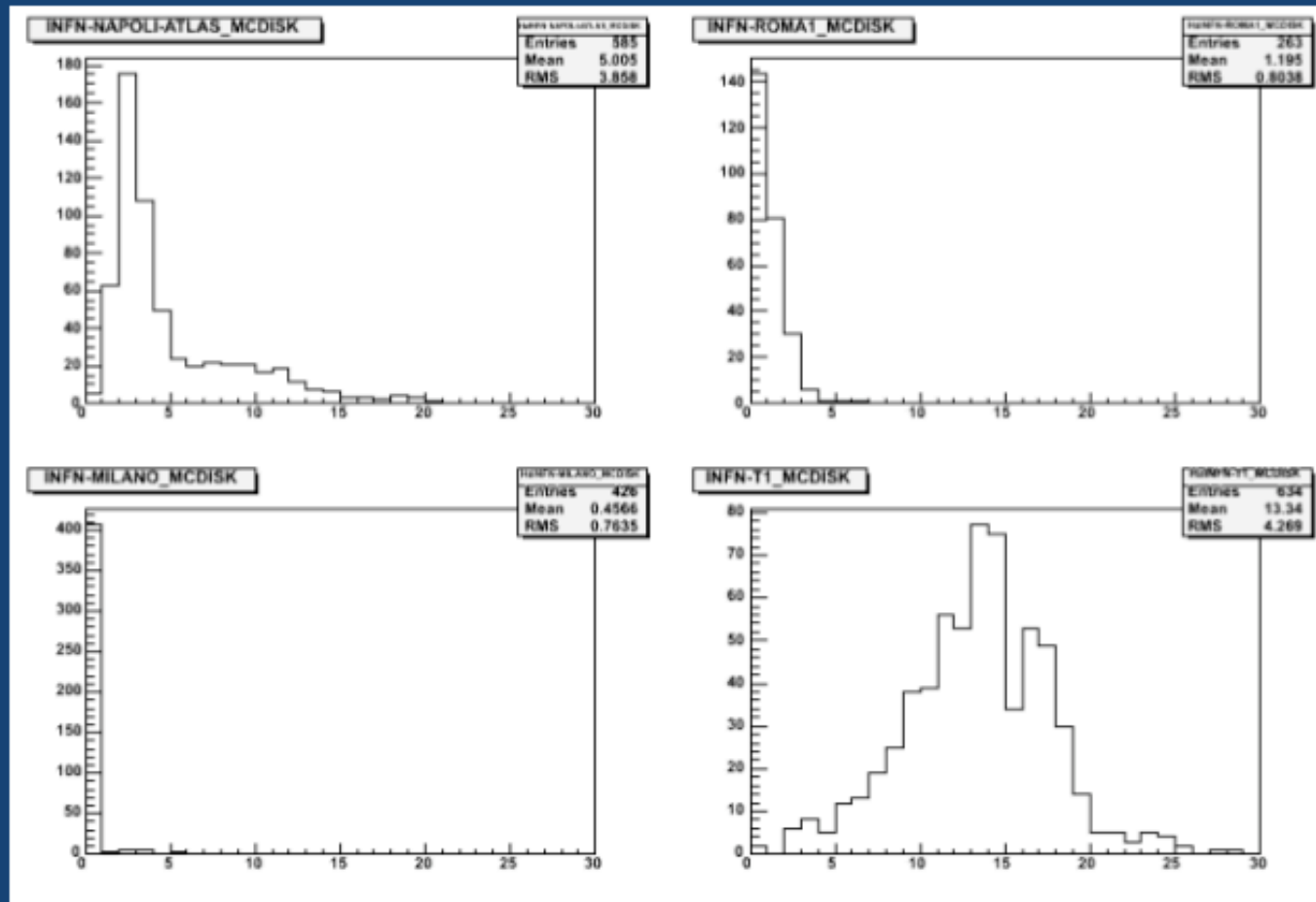
## *“Prima” Efficiency*

Overall Efficiency (excluding expired proxies):

Site	Status				C/(C+F)
	TOTAL	COMPL	FAILED	RUNNING	
INFN-MILANO_MCDISK	442	320	107	15	74.90%
INFN-NAPOLI-ATLAS_MCDISK	588	391	196	1	66.60%
INFN-ROMA1_MCDISK	360	257	7	96	97.30%
INFN-T1_MCDISK	661	631	30	0	95.50%

CNAF was good. Milano/Napoli saw ~30% failures.  
Roma was quite slow, but not failing many.

# “Prima” Events/Second

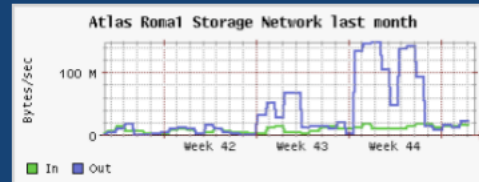


CNAF: 13 Hz, NAPOLI: 5 Hz, ROMA: 1 Hz, MILANO: 0.5 Hz



# Leçons et limitations

- Milano, Napoli, Roma all had saturated 1Gb networks:



- Assuming 0.2MB per event, a 1Gb network can stream at 640 Hz. With 200 CPUs, we expect, and observed  $\sim 3$  Hz.
- Other notes:
  - Milano found a faulty switch (replaced afterwards).
  - Napoli has a 10Gb switch connected to only 2 racks.
  - Roma is physically inaccessible.

## □ Limitations du réseau

### □ 'Trop' de CPUs

## □ Nécessité d'avoir un contact interne sur site

### □ Pour monitoring et diagnostique

Soient  $\sim 300K$   
evts par jour  
C'est peu...  
Data taking  
200Hz



## Results

Site	Submitted	Running	Completed	Failed	Total
CYFRONET-LCG2.MCDISK	0	0	391	9	400
DESY-HH.MCDISK	0	0	296	4	300
DESY-ZN.MCDISK	0	0	380	20	400
HEPHY-UIBK.MCDISK	0	0	185	100	285
LRZ-LMU.MCDISK	0	0	369	31	400
PRAGUELCG2.MCDISK	0	0	164	136	300
WUPPERTALPROD.MCDISK	0	0	213	87	300
UNI-FREIBURG.MCDISK	0	0	17	14	31

Very low efficiency for analysis job on some sites

- For all sites more than 50% of the jobs succeeded. For 4 sites even better than 90%
- Many (>100) errors : Error You cannot open a ROOT file in mode READ if it does not exists., problem seems to be due to a failing lcg-gt to determine the file TURL.
- Other errors need to be investigated.

Meme conclusions pour nuage DE

Limitation of internal bandwidth and site configuration

## CPU/Walltime

- 2 sites (CYFRONET, PRAGUE) have a very bad CPU/Walltime due to network limitation : only 1Gbps links to the pools.
- Best performance for DESY-HH which is the biggest site and has 75% of all AOD (better spread of the data on the pools ?)
- For other sites the data are stored only on a few pools : e.g. : 1 in DESY-ZN, 3 in LRZ-LMU. Some saturation observed (see next slide)

# Xrootd in ATLAS

46

- SLAC :
  - ▣ data storage,
  - ▣ SRM interfaced
- BNL :
  - ▣ data storage, No SRM interface
  - ▣ PROOF
    - <http://www.usatlas.bnl.gov/twiki/bin/view/AtlasSoftware/ProofTestBed.html>
- Wisconsin :
  - ▣ PROOF + Condor
- CC-Lyon :
  - ▣ data storage, No SRM interface

Découplage stockage /  
distribution des données :  
fonctions propres du T1  
Accès intensif : Ferme  
d'Analyse



# Actions pour 2009 (et fin 2008)

47

ATLAS Software and Computing workshop @CERN November 3-7 2008

## Timing

- November-December
  - preparation for re-processing
  - scale-up MC and Analysis FT
  - Super throughput challenge
  - Panda installation at CERN
- January-March
  - re-processing
  - further scale-up MC and Analysis FT
  - Panda running from CERN now
- April
  - data clean-up
  - preparations for data
- May-June
  - cosmics and detector commissioning
- July – December
  - collisions

- Stagein does not work on most dCache T1s
- Lyon does not meet requirements
- Scalability issues of ORACLE setup at some T1s observed

Analysis Functional test  
cad. Analysis Challenge en continue

T1 - T1 Data exchange

De la difficulté de faire des prédictions  
L'optimisme biaise toujours...

11/06/08

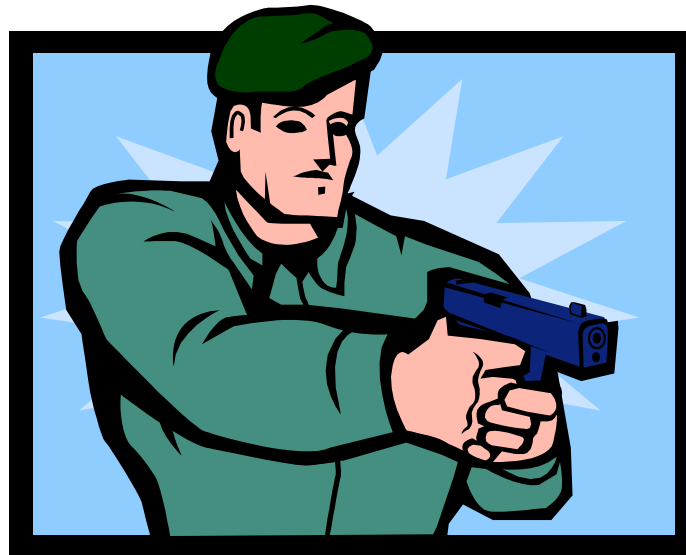


# J'ose des conclusions

48

- Aurions nous été prêts si le LHC avait délivré des données
  - ▣ NON?
- Nous avons progressés moins vite que la majorité des autres nuages d'ATLAS
  - ▣ Malgré un bon départ en 2006
- Le T1 est un soucis
  - ▣ Stockage & distributions données
  - ▣ Nb de FTEs directement affectés à l'opération du T1 ( par ex : CE + dCache + HPSS + FTS + LFC < 6 FTE)
- Les problèmes n'ont peut être pas encore atteints les T2s
  - ▣ Pas d'analyse...







# Panda Job States

50

□ 10 values in Panda describing different possible states of jobs

- **defined** : job-record inserted in PandaDB
- **assigned** : dispatchDBlock is subscribed to site
- **waiting** : input files are not ready
- **Activated** : waiting for pilot requests
- **sent** : sent to a worker node
- **running** : running on a worker node
- **holding** : adding output files to DQ2 datasets
- **transferring** : output files are moving from T2 to T1
- **finished** : completed successfully
- **failed** : failed due to errors

□ Normal sequence of job-states:

- **defined -> assigned -> activated -> sent -> running -> holding -> transferring -> finished/failed**



# Data Volume and Rates

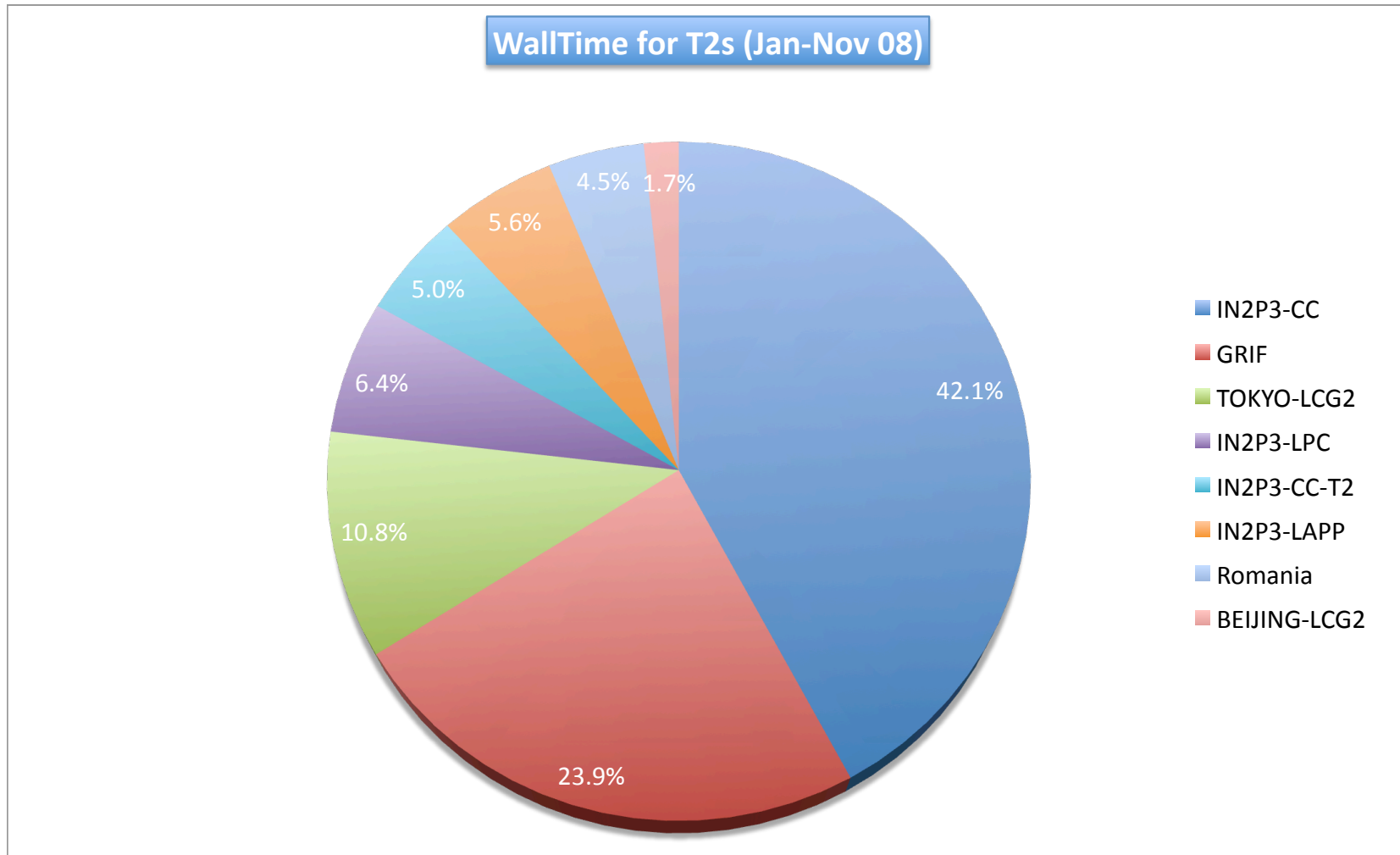
51

Trigger Rate	200 Hz	
RAW	1.6 MB/evt	~28 TB/day
ESD	1.0 MB/evt	~17 TB/day
AOD	0.2 MB/evt	~3.5 TB/day



# CPU consommation par site du nuage

52





# PanDA Server

