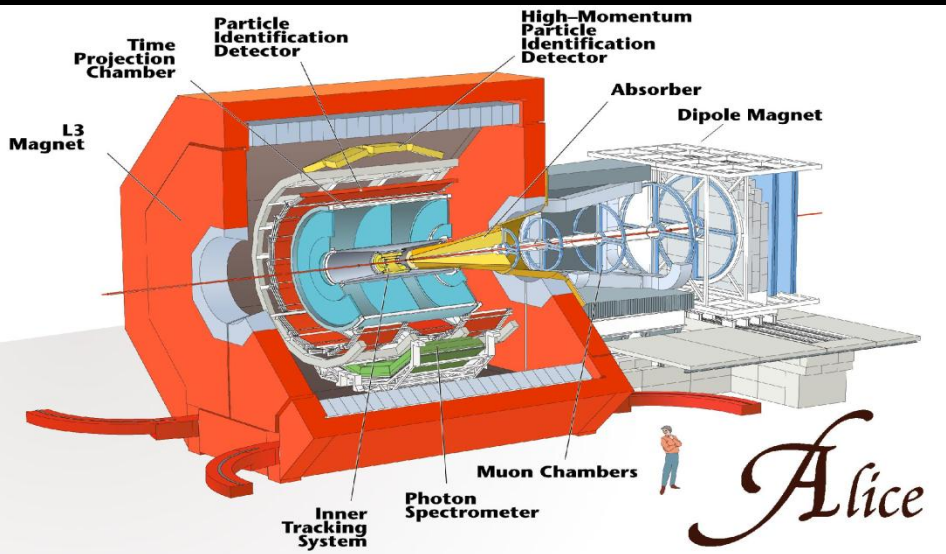




# The ALICE Grid operations

T1-T2 Workshop, KIT  
24/02/2012



## ALICE Collaboration

- ~ 1/2 ATLAS, CMS, ~ 2x LHCb
- ~1100 people
- 30 countries,
- 80 Institutes

Total weight	10,000t
Overall diameter	16.00m
Overall length	25m
Magnetic Field	0.4Tesla

8 kHz (160 GB/sec)

level 0 - special hardware

200 Hz (4 GB/sec)

level 1 - embedded processors

30 Hz (2.5 GB/sec)

level 2 - PCs

30 Hz

(4 GB/sec)

data recording &  
offline analysis

## The ALICE collaboration & detector

# Data volumes

- RAW data – 2.5 PB/year
  - Two distinct periods – p+p (~7.5 months) and Pb+Pb (~40 days)
- Reconstructed and simulated data
  - 1.5PB – first level RAW filtering (ESDs)
  - 200TB – second level RAW filtering (AODs)
  - 1PB of simulated data
- User generated data ~500TB
- **Total ~5 PB of data per year (without replicas)**
  - Replication 2x RAW, 3x ESD/AODs, 2x user files

# Processing

- RAW data reconstruction ~10K CPU cores
- MC processing ~15K CPU cores
- User analysis ~7K CPU cores (450 distinct users)
- ~40Mio jobs per year
  - ~ 1.3 job completed every second
    - $\frac{1}{2}$  production,  $\frac{1}{2}$  user jobs
- 200 Mio files per year

# How do we go about that?



72 active computing sites



# The Grid structure

- T0 – CERN
  - Stores all RAW data from the 4 LHC experiments (custodial)
  - First pass reconstruction, MC, user analysis, interactive analysis, detector calibration
- T1s – Large regional computing centres
  - Stores part of the RAW, ESDs, AODs (custodial)
  - Second pass reconstruction, MC, user analysis and interactive analysis facilities (AFs)
- T2s – computing centres of universities/labs
  - MC production, user analysis

# The Grid structure (2)

- Various flavours of site middleware
  - gLite, ARC, AliEn direct to batch
- Storage is xrootd-enabled
  - CASTOR, dCache, DPM, xrootd native (majority of the SEs)
- All is interfaced through the AliEn central services: single file catalogue and central task queue for job management
- The users are fully 'protected' from the Grid plumbing
  - In this sense the ALICE Grid is a cloud
  - Fully transparent data access (when possible) would make it a 'super cloud'

# Some history

- Working prototype in 2002
- The Vision from the very beginning
  - Single interface to distributed computing for all ALICE physicists
  - File catalogue, job submission and control, application software management, end user analysis
  - And this is....

**AliEn – Alice Environment**

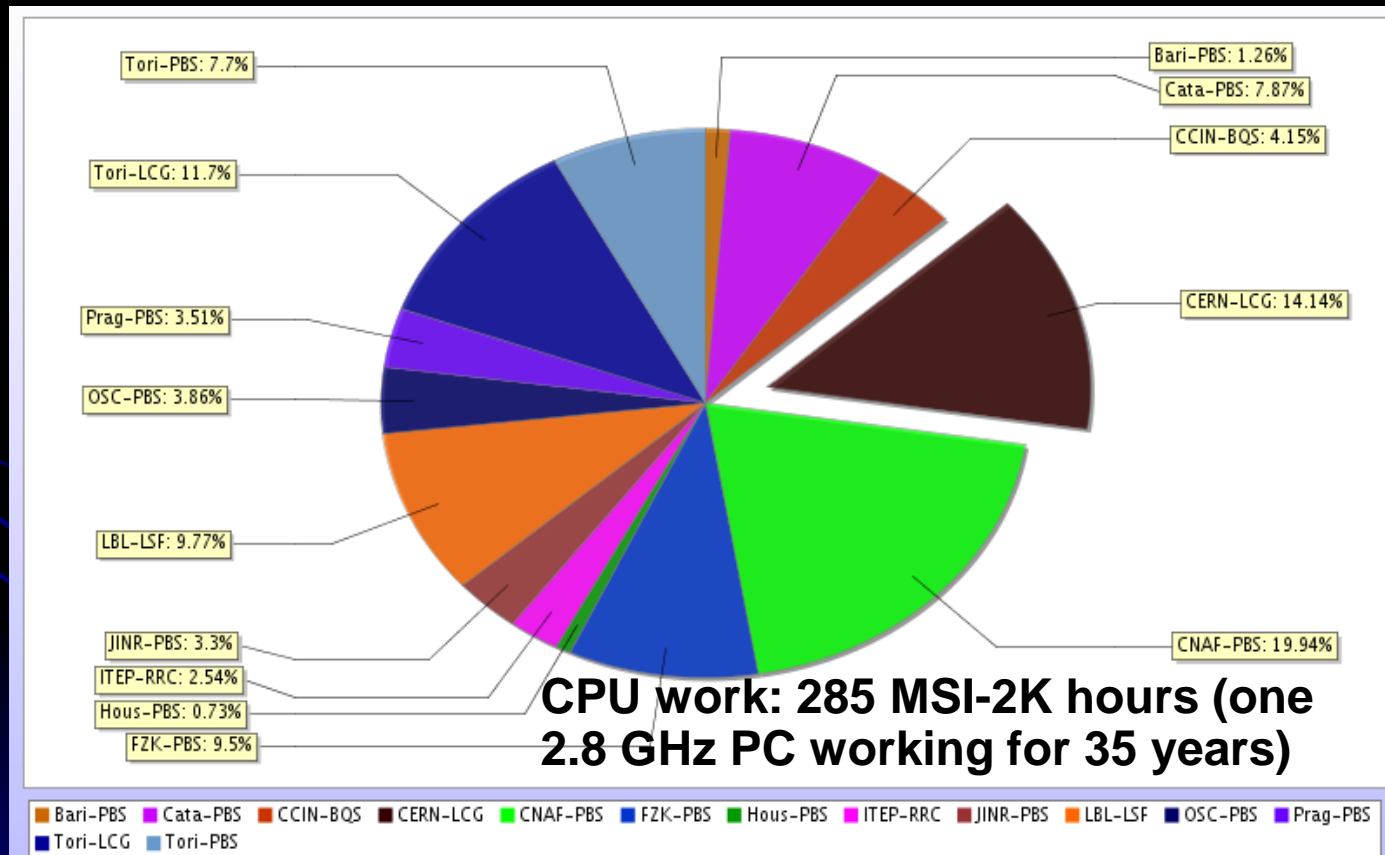


# Grid use in the first years

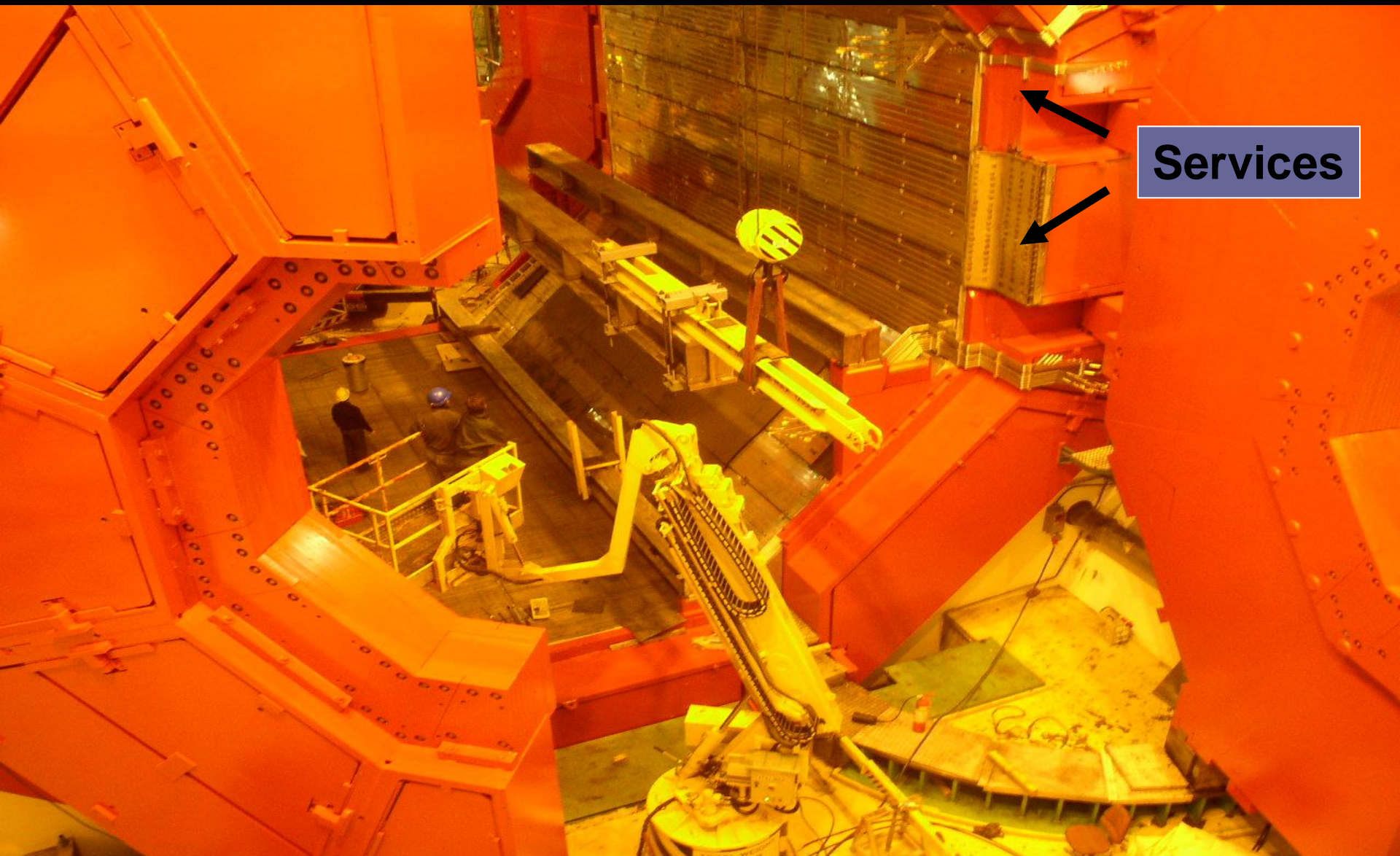
- First MC productions...
- Vertical Grid – interfaced down to local batch system level and any type of local storage the site provides (capability retained and refined up to today)
- Few hundred CPUs at 13 sites
- Very ambitious goals – validation of the entire ALICE Computing model

# First distributed production

- First report of yearly ALICE data challenge: Sep. 2004

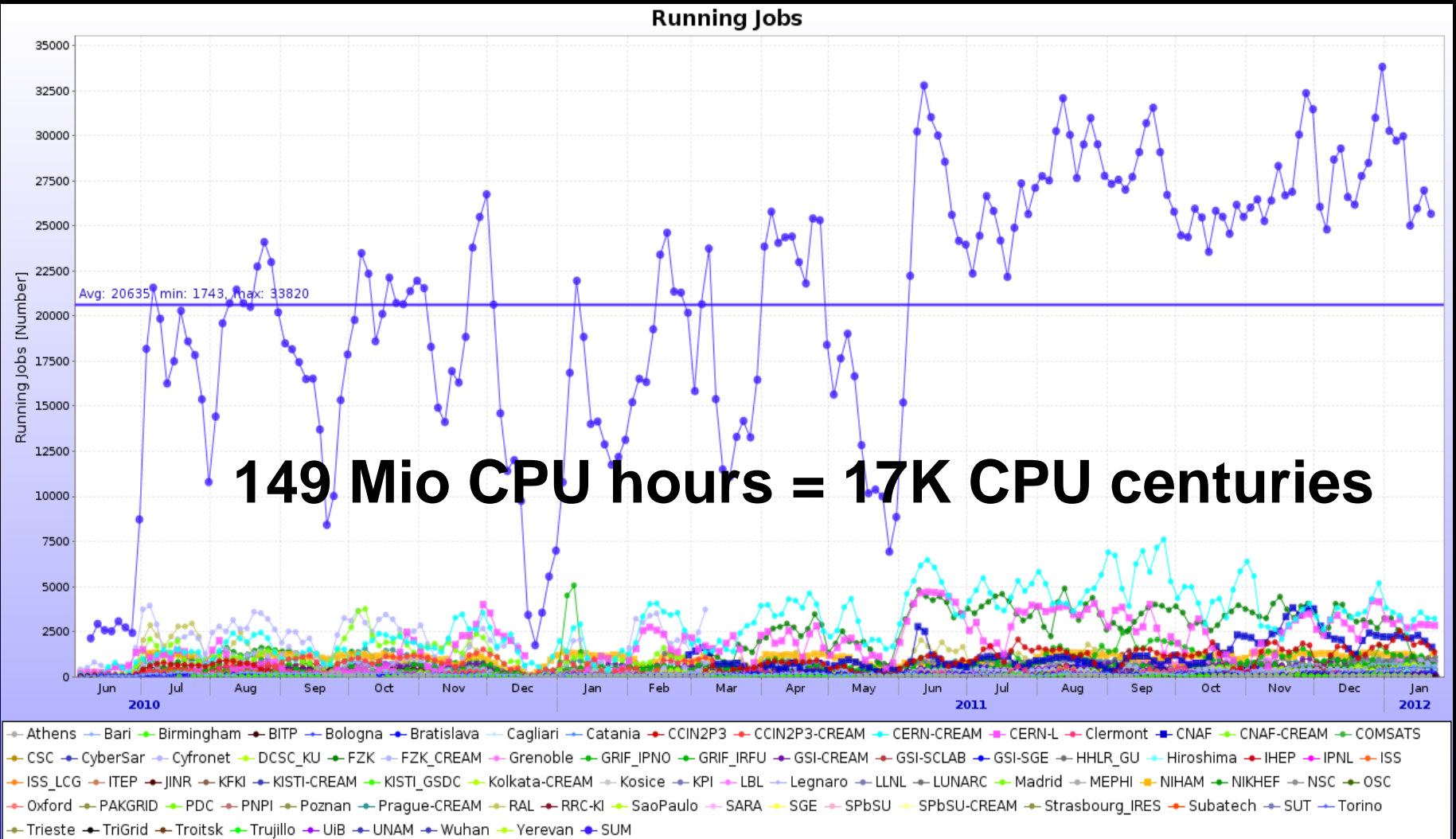


# The experiment building begins

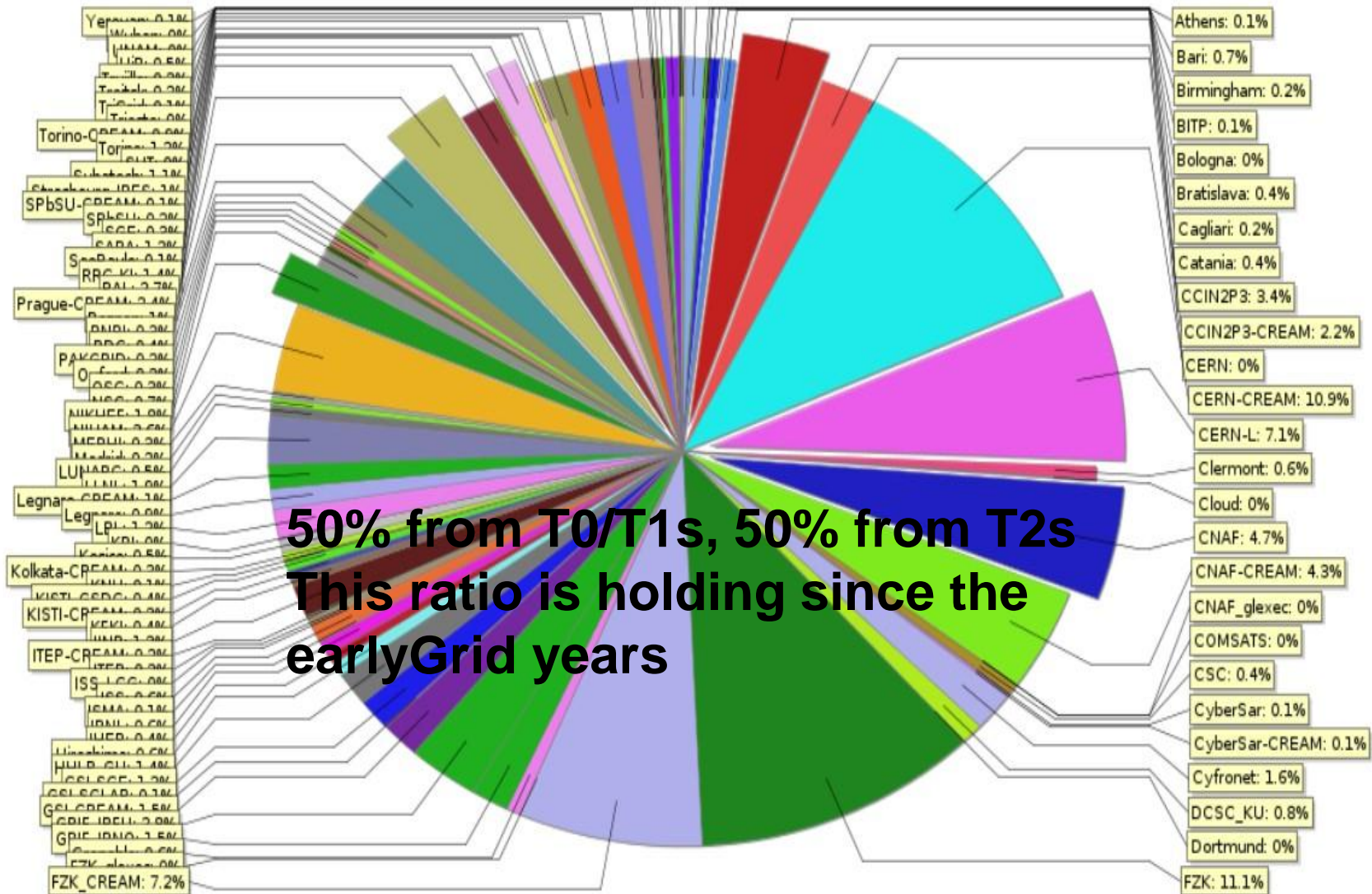


Services

# The Grid from 2010 onward



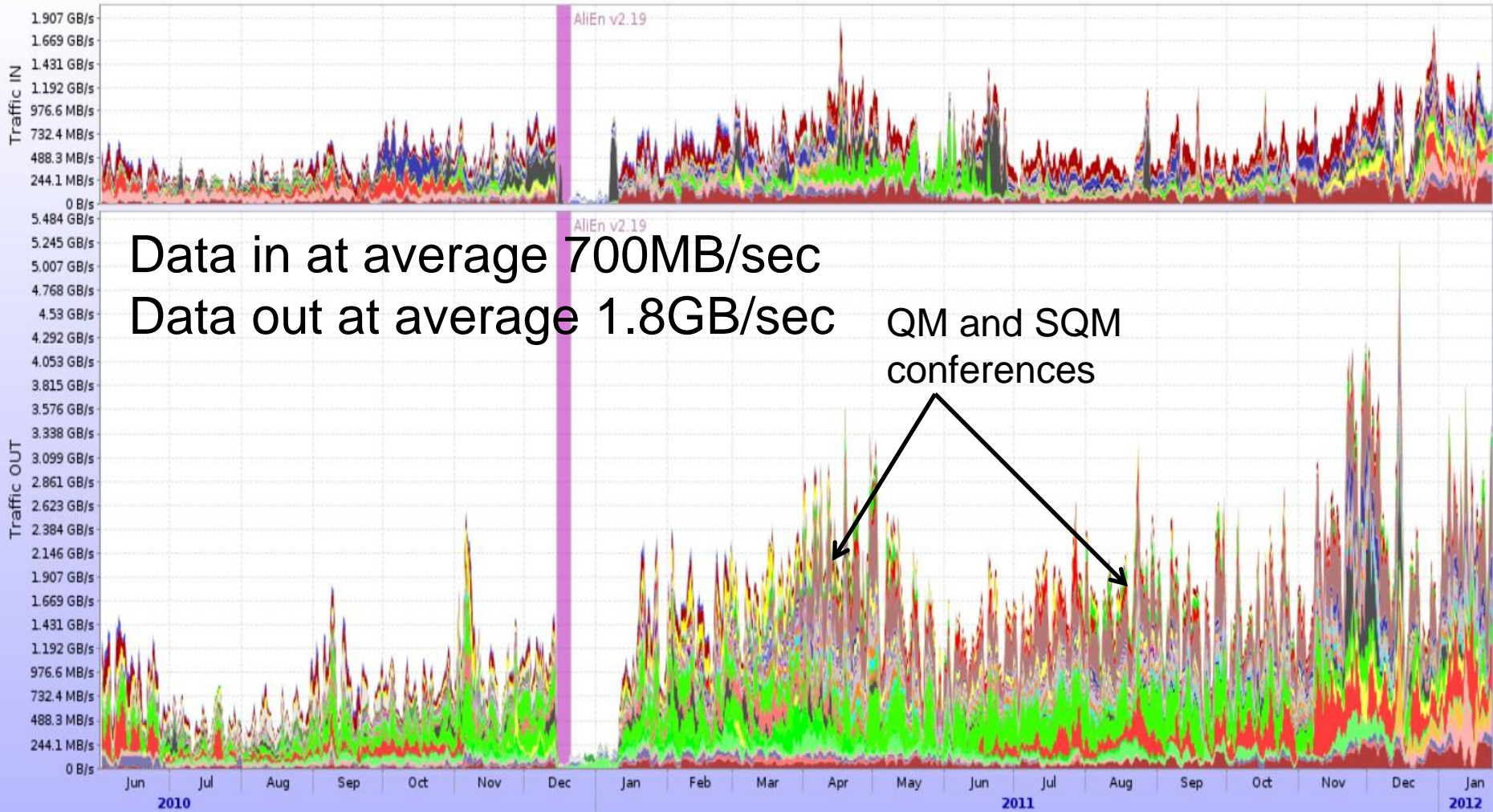
# Centres contribution



**50% from T0/T1s, 50% from T2s**  
**This ratio is holding since the**  
**earlyGrid years**

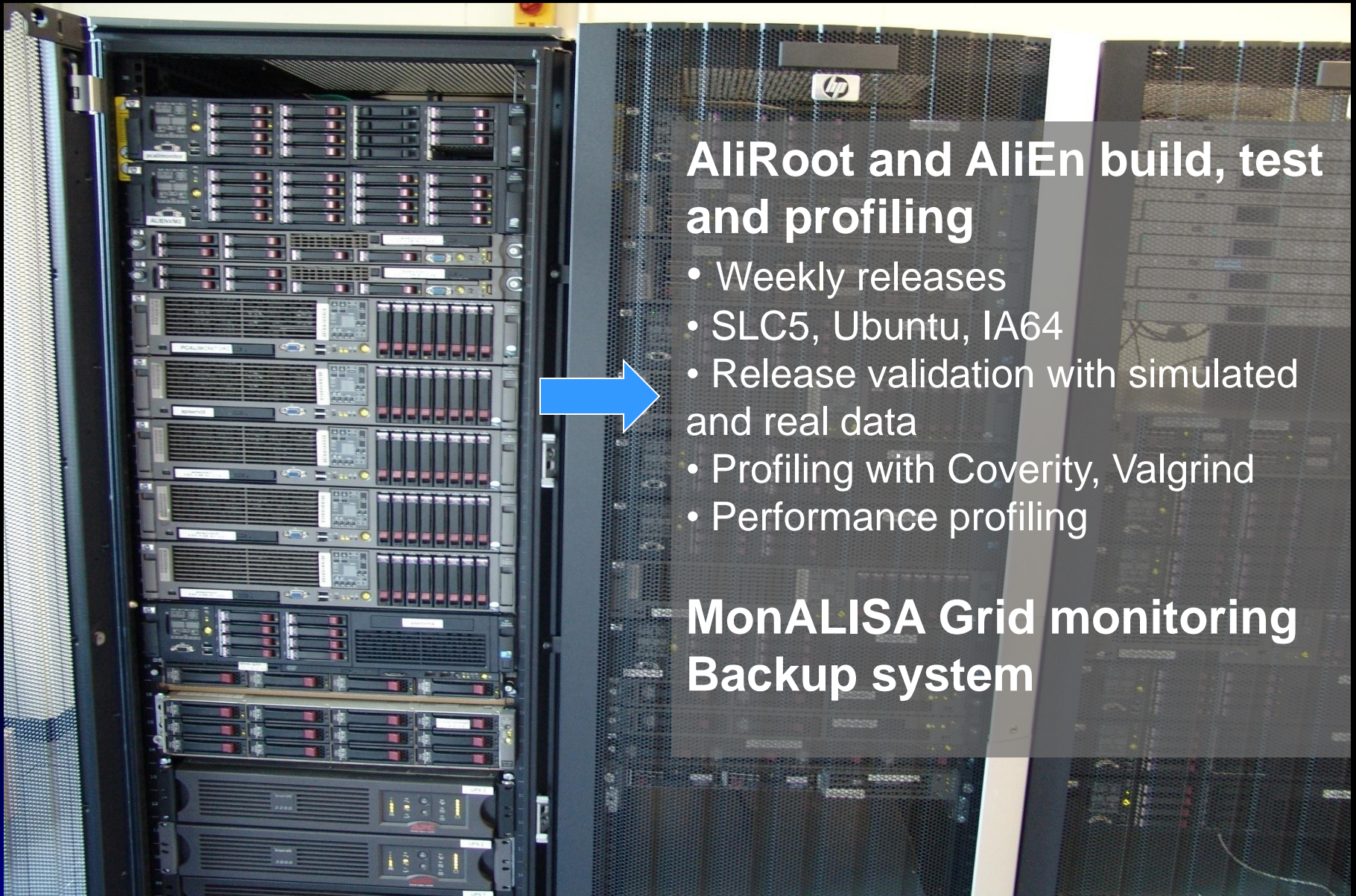
# Data r/w

Aggregated network traffic per SE



# How is all this managed

- AliEn central services (see Pablo's presentation)
- Site services and support
- Monitoring (see Iosif's presentation)
- PROOF enabled AFs (several presentations)
- Organization of meeting/contact lists
- Challenges

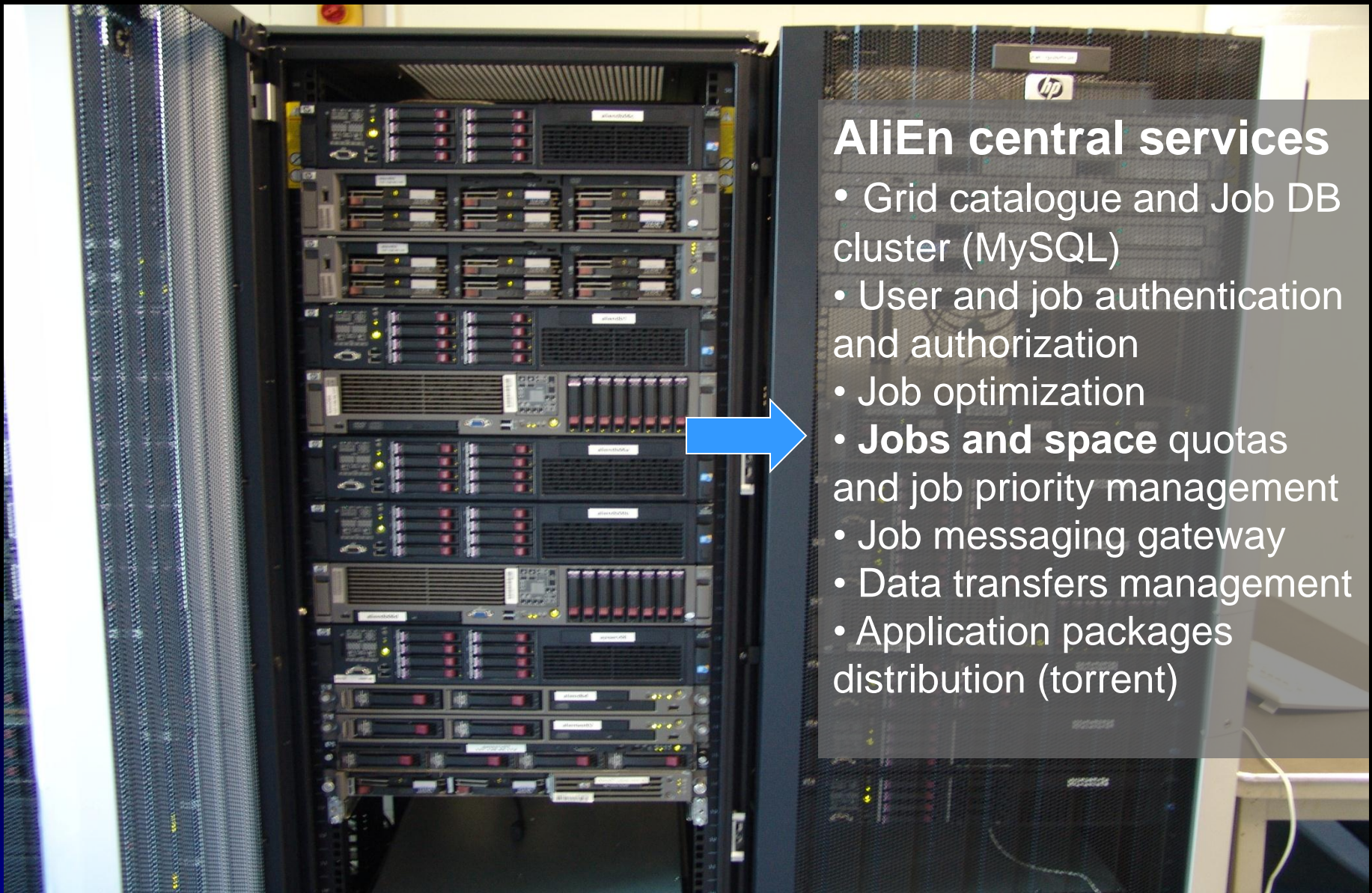


## AliRoot and AliEn build, test and profiling

- Weekly releases
- SLC5, Ubuntu, IA64
- Release validation with simulated and real data
- Profiling with Coverity, Valgrind
- Performance profiling

## MonALISA Grid monitoring Backup system





## AliEn central services

- Grid catalogue and Job DB cluster (MySQL)
- User and job authentication and authorization
- Job optimization
- **Jobs and space** quotas and job priority management
- Job messaging gateway
- Data transfers management
- Application packages distribution (torrent)

AliEn and AliRoot  
MAC OSX build and  
test system

ALICE web server  
(<http://aliweb.cern.ch>)

AliEn job and user  
connection API  
services

Central services



# Central services

- Up 24/7, 365 days... no other service in the experiment comes close
  - Critical for data taking
  - Critical for Grid processing in general
- General consideration
  - All services are running multiple instances and are load-balanced
  - Good hardware, lights-off operation
- Still, they are a 'single point of failure'
  - As demonstrated during extended power cuts
- Recovery takes 2 hours from cold start
  - Over the years we have mastered the operation
  - Having a truly 'distributed' central services is possible, but costly...

# Site services (1)

- VO-box component
  - 1gLite (proxy renewal) + 5 AliEn (CMreport, ClusterMonitor, PackMan, CE, MonALISA)
  - AliEn services are monitored and send alarms (if site expert is subscribed to it) in case of trouble
  - All services can fail with almost equal probability, auto-restart does not always work
  - Human intervention is sometimes necessary
  - In general stable...

# Site services (2)

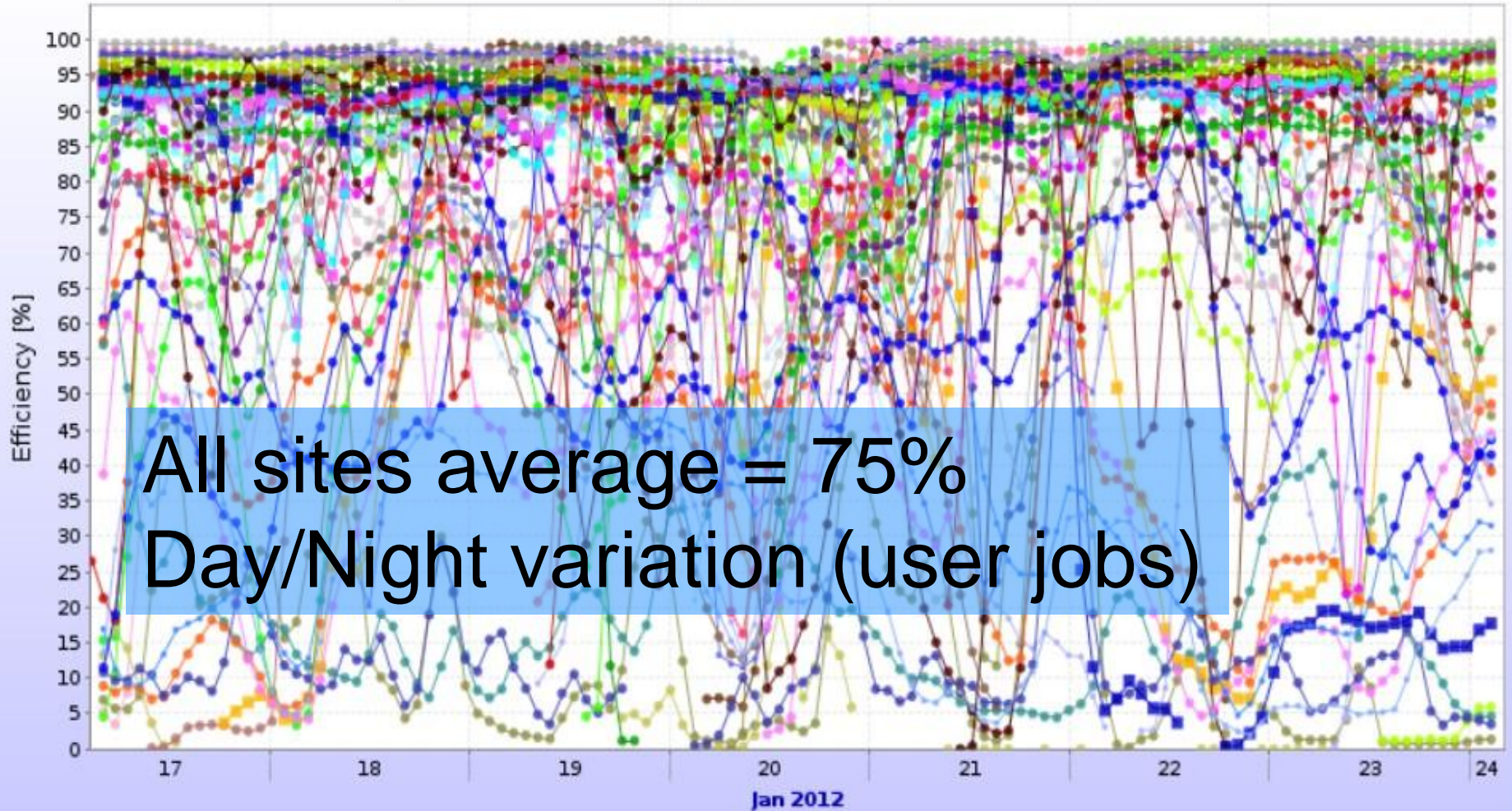
- Site SE
  - Functional tests through central 'ping' test
  - Checks the 'general' SE working status, but not all servers
  - **The SE is by far the most critical component to watch (single point of failure for every site)**
  - Also instrumented with alarms, please subscribe!
  - Individual server alarms would be a bonus, but presently we don't have them (most of the sites do not have local monitoring too)
  - See Iosif's talk for further info...

# More on storage

- After 2 years of RAW/MC/User activities, the disk storage elements are getting quite full
  - Negative effect on sites performance— less ‘useful’ work done, lower efficiency due to remote data access (all write/some read)
- Ideally, all SEs would have the ‘same’ occupancy, however
  - Not all SEs came to ‘life’ at the same time
  - CPU/SE ratio is not optimal at many sites
  - The auto-discovery mode helps, but locality plays a role too...inter-site networks are not yet allowing for fully transparent SE use

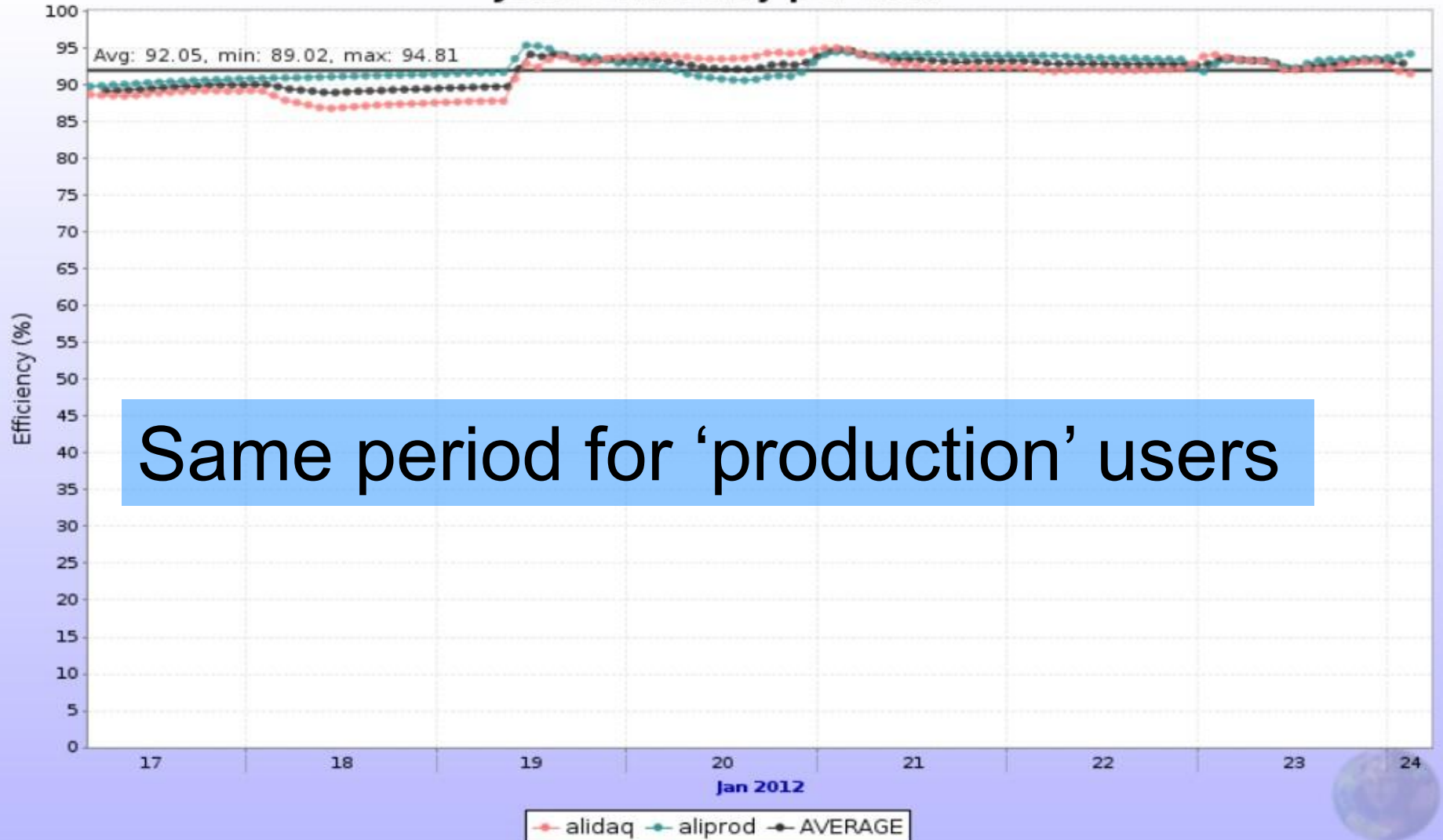
# A word on efficiency

Jobs efficiency (cpu time / wall time)



# A word on efficiency (2)

Jobs' efficiency per user



Same period for 'production' users



## More on storage (2)

- Removal of data ongoing in 3 directions
  - Reduction of replicas for older RAW data production cycles
  - Reduction of replicas for less-used and less-important MC productions
  - Stricter user quotas
- In all cases, the battle is tough and the data generation far outpaces any removal
- Full provision of pledged resources will balance back the sites

# Organization

- Single mailing list (for most communications)
- Direct contact with regional experts/site experts (both ways) works well and is fast
- Weekly AF/TF meetings
  - Proposal to make them bi-weekly (or monthly)
    - Grid is generally stable – but longer and with specific agenda for every meeting
  - Only makes sense if experts participate

# Challenges

- Upgrade AliEn to v.2-20
  - Discussion with ALICE Physics Board
  - Reconstruction (Pass2) of 2011 Pb data must finish
- Improve the situation with the site SEs
  - Stability and occupancy
- Improve the efficiency of user jobs
  - Partially the above will help
  - For more details (roadmap) see Andrei's presentation

# Conclusion

- The Grid has reached excellent level of maturity for most application
  - Work is needed to be able to handle efficiently user jobs, as these will continue to be a major part of the resources utilization
  - In this direction – modification of job assignment and handling in AliEn
- Mastering the storage is still a challenge
  - Stability of every individual server is a must
  - More extended monitoring and new xrootd version will help

# Conclusions (2)

- General improvements in the network (LHCOPN) will further dissolve the tier boundaries
  - Meanwhile, better routing especially for the T2s would help
  - Network is a sufficient topic for a separate workshop
- As a whole, the Grid is fulfilling its role and enables the ALICE physicists to carry out immense computational tasks in a distributed computing environment