



# ALICE Data Processing



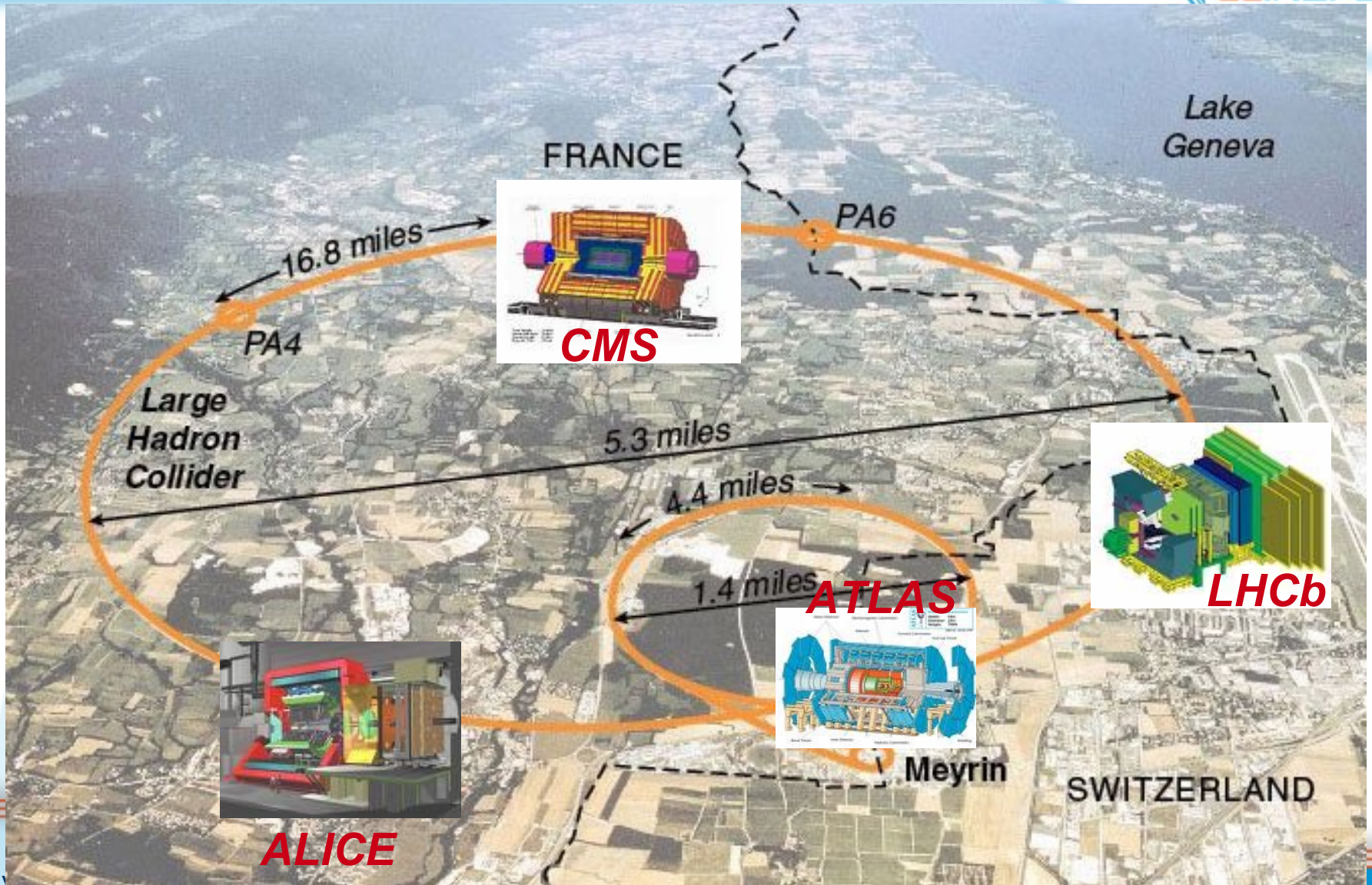
# Outline



- ALICE experiment
- Data acquisition, processing & transfer
- AliEn, job submission, workload management...
- What should happen at CC-IN2P3



# The ALICE experiment at LHC





# The ALICE experiment

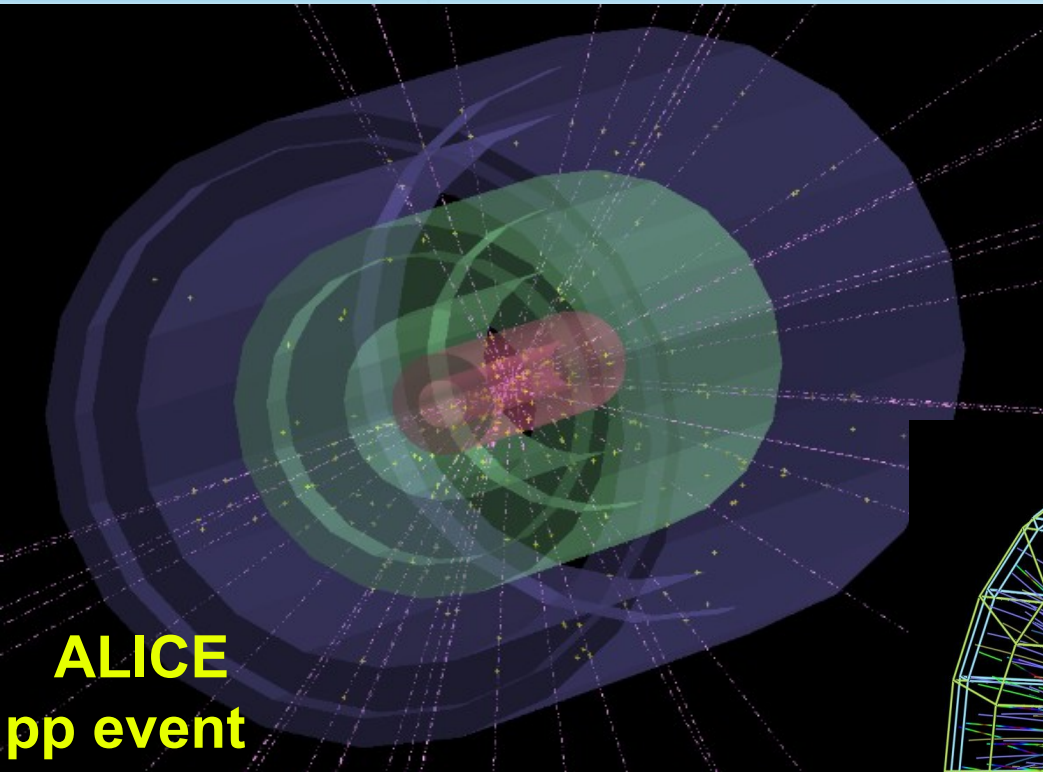


- Only experiment at LHC originally dedicated to the study of the Quark-Gluon Plasma
  - ◆ Heavy Ion collisions ! (Pb+Pb)
    - Best way to create the QGP : very large energy density available
  - ◆ Proton-proton collisions
    - Reference for heavy ions and genuine pp physics
  - ◆ Most of the time, LHC will run pp collisions (1 month of HI/year)
- ALICE computing model has been designed according to the needs of heavy ion collisions
  - ◆ Large amount of data stored per Pb+Pb collision
    - Expected mean rate: 1.25 GB/s data produced
  - ◆ Very long computing time (simulation, reconstruction, analysis)
  - ◆ Much storage space

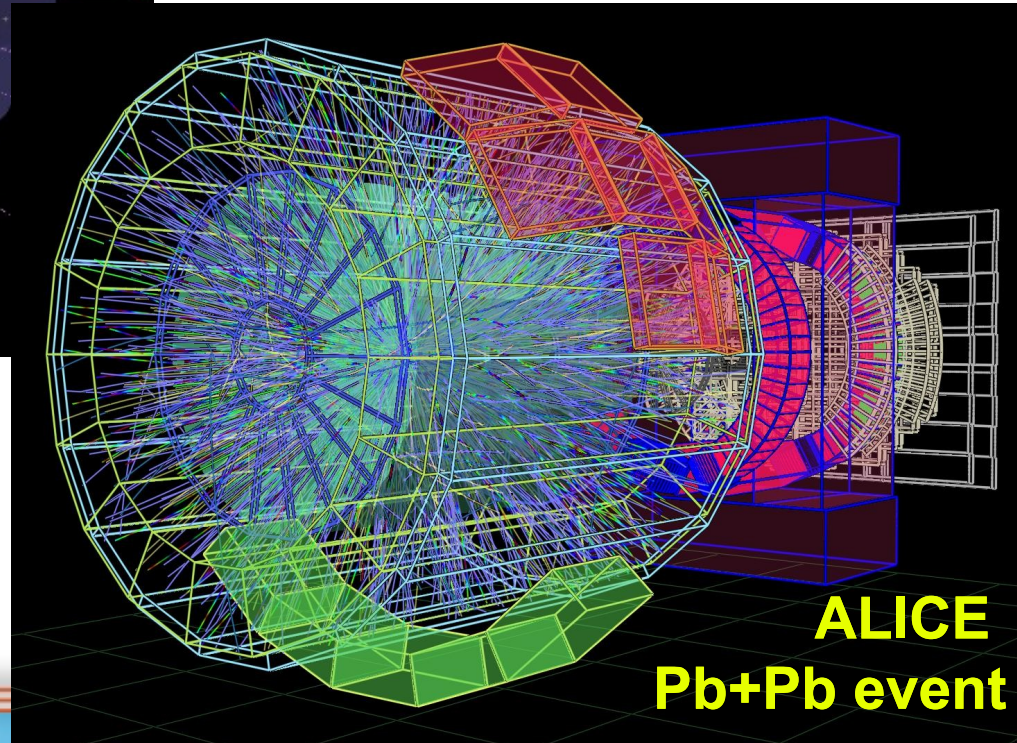




# ALICE event snapshot



**ALICE  
pp event**



**ALICE  
Pb+Pb event**



# Computing Model – pp

(Permanent Data Storage)

PDS



(AliEn File Catalogue)

AliEn  
FC

CERN T0



ODB

First pass reconstruction

CAF Analysis (PROOF)  
External Copy

RAW & Calibration parameters

(Online Data Buffer)

(Event Summary  
Data)

ESDs

Ext T1s



Pass 1 & 2  
reco

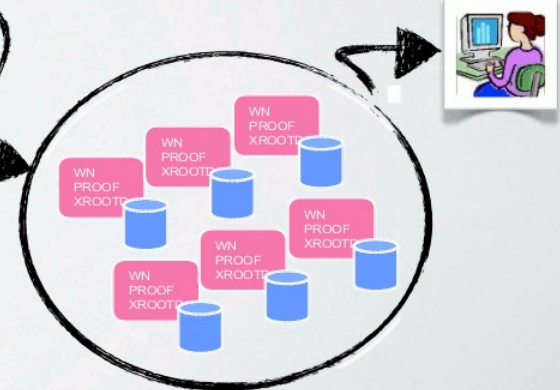
Analysis  
Train

Ext T2s

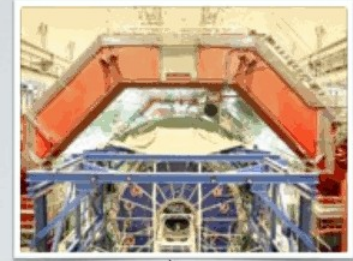


MC

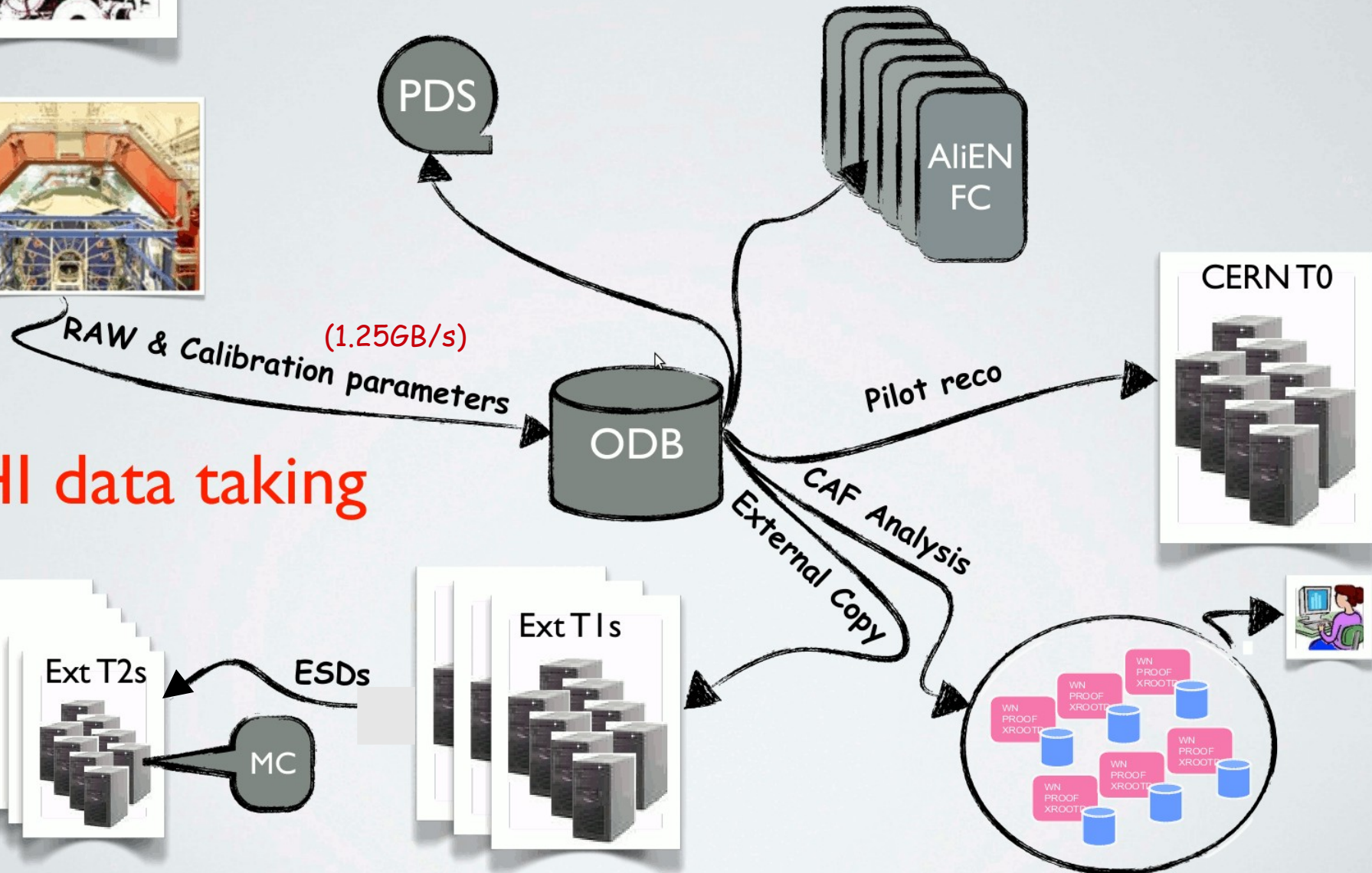
End users  
analysis



# Computing Model – AA



HI data taking

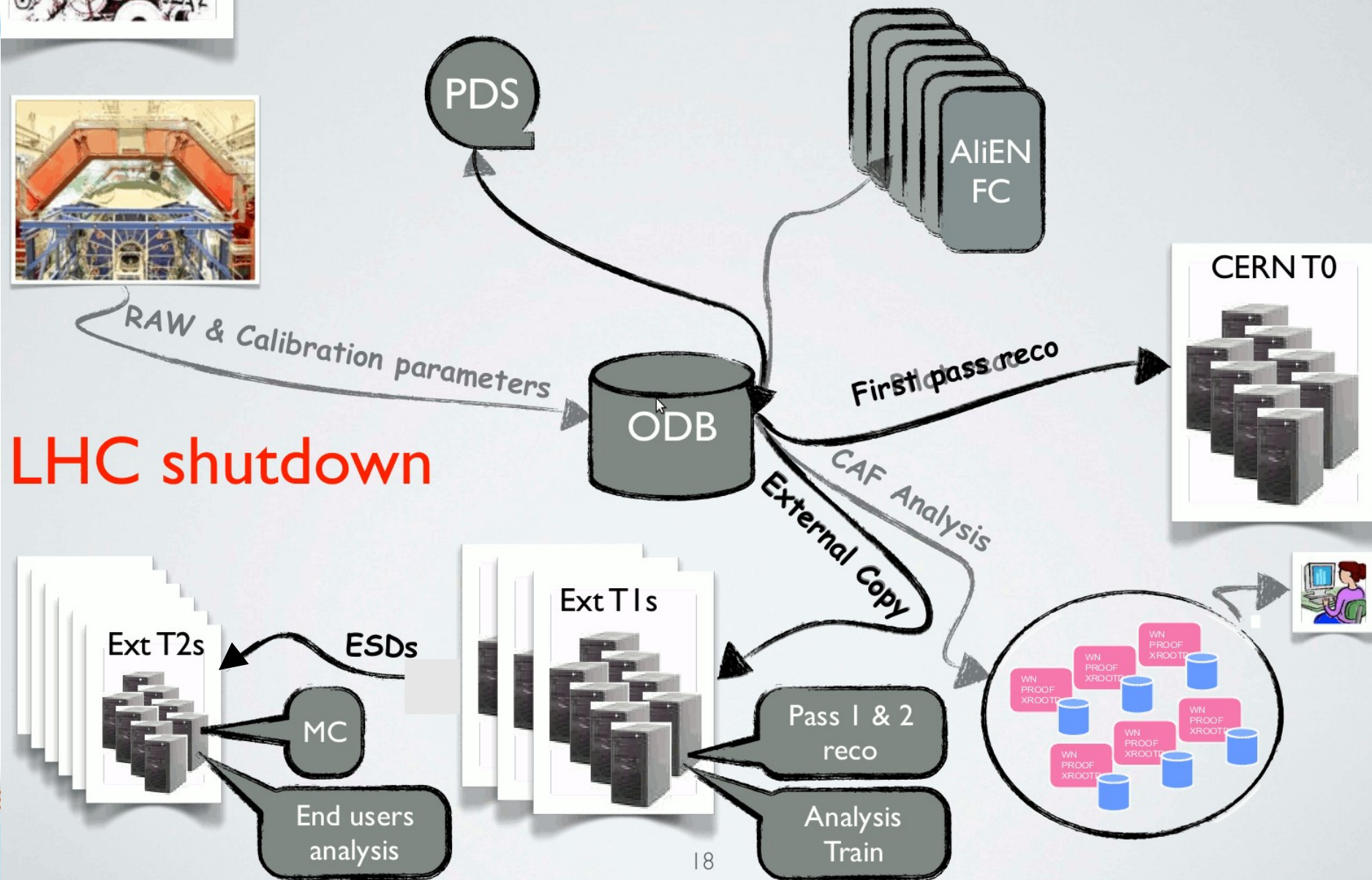




# Computing Model – AA



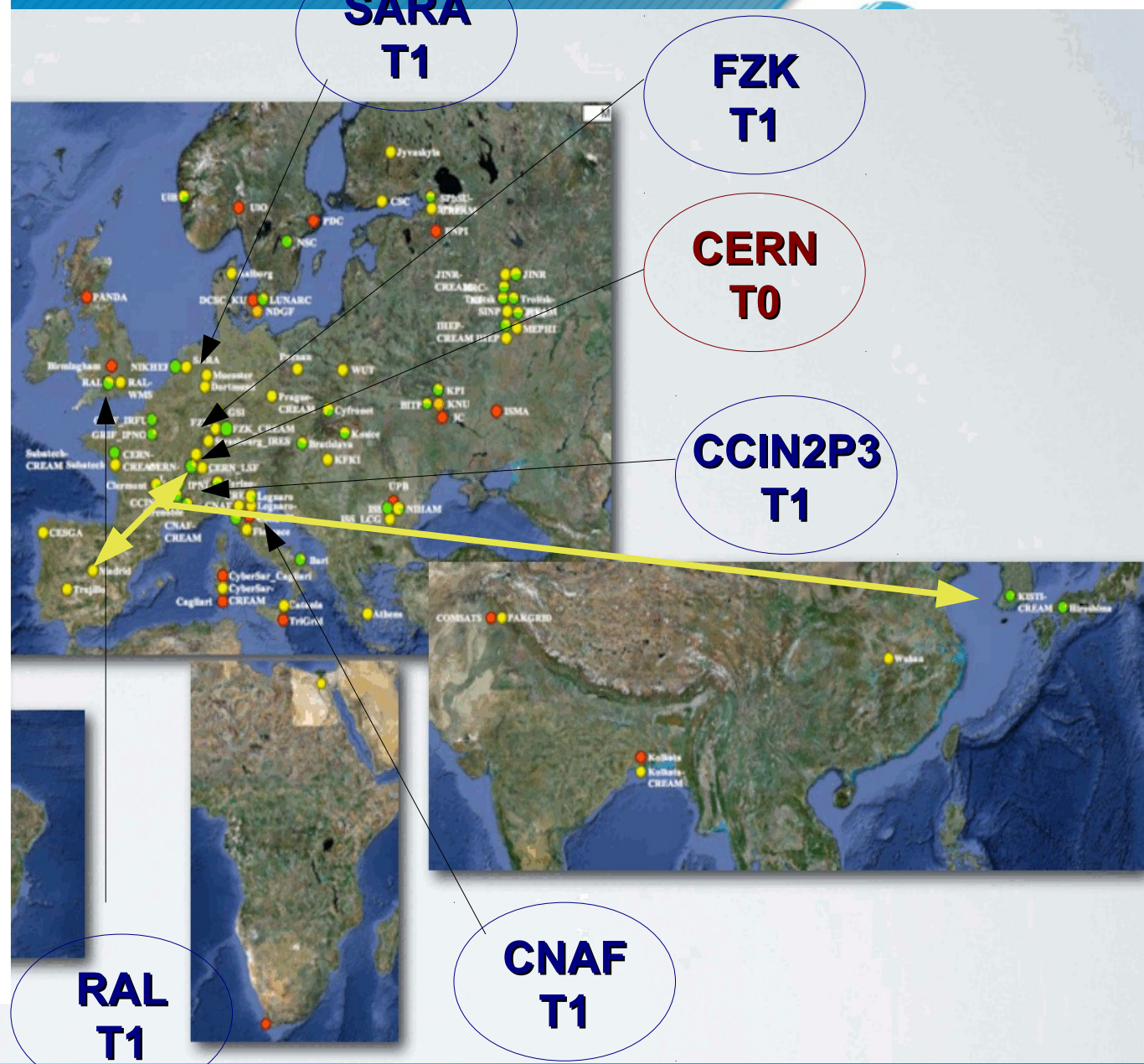
LHC shutdown





<b>Tier 1</b>	<b>Tier2</b>
<b>CCIN 2 P3</b>	<b>French Tier2</b>
	<b>Sejong (Korea)</b>
	<b>Lyon Tier2</b>
	<b>Madrid (Spain)</b>
<b>CERN</b>	<b>Cape Town (South Africa)</b>
	<b>Kolkata (India)</b>
	<b>Tier2 Federation (Romania)</b>
	<b>RMKI (Hungary) <sup>2</sup></b>
	<b>Athens (Greece) <sup>2</sup></b>
	<b>Slovakia</b>
	<b>Tier2 Federation (Poland)</b>
	<b>Wuhan (china)</b>
<b>FZK</b>	<b>FZU (Czech Republic)</b>
	<b>RDIG (Russia)</b>
	<b>GSI (Germany)</b>
	<b>Muenster (Germany)</b>
<b>CNAF</b>	<b>Tier2 Federation (Italy)</b>
<b>RAL</b>	<b>Tier2 Federation (UK)</b>
	<b>Birmingham <sup>2</sup></b>
<b>NIKHEF</b>	<b>SARA<sup>2</sup></b>
<b>PDSF<sup>1</sup></b>	<b>LLNL (USA)</b>
	<b>OSC (USA)</b>
	<b>Houston</b>

# E GRID





# The data transfers

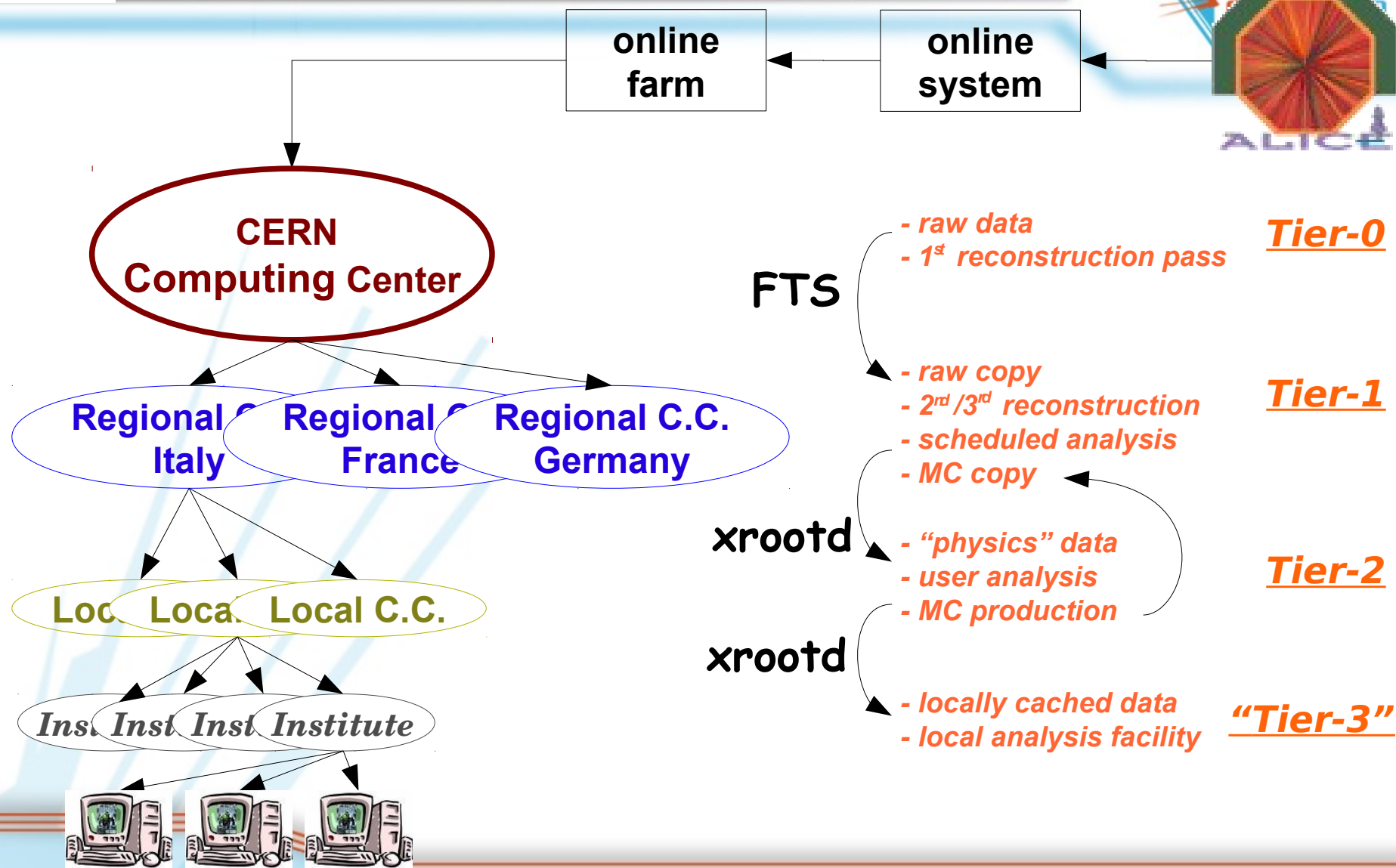


- All scheduled transfers use FTD (~AliEn's FTS)
  - ◆ Transfer Queue
  - ◆ Triggered by alien 'mirror'
- $T0 \rightarrow T1$  and  $T1 \rightarrow T1$  use LCG's FTS
  - ◆ Defined channels
  - ◆ Data go in and out the SE's via SRM interface
- $T1 \rightarrow T2$  and  $T2 \rightarrow T2$  use xrootd
  - ◆ No predefined channel
  - ◆ Data go in and out the SE's via xrootd server





# Data transfers and analysis (in brief)

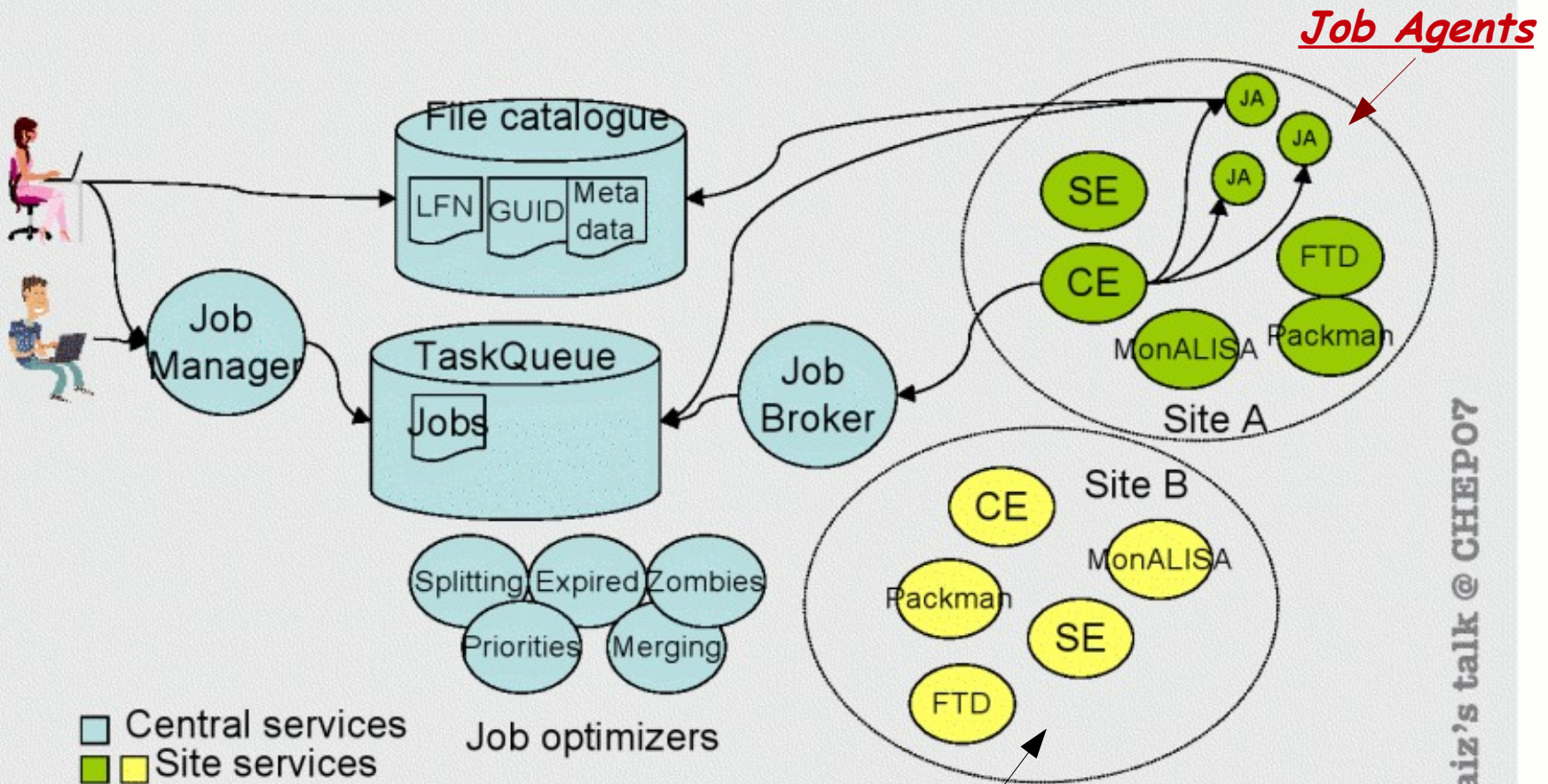




- GRID middleware for ALICE
- AliEn provides
  - Site services (CE/SE interfaces, package manager...) on VObox
  - File catalogue
    - All ALICE data (simulation+reconstruction)
    - User data
  - Client services
    - Catalogue browsing
    - Job submission
    - Job monitoring
    - Data transfers : between SE's + local→SE + SE→local



# JOB EXECUTION MODEL



Pablo Saiz's talk @ CHEP07



# AliEn (user view)

```
rvernet@ccvcre:/$ aliensh
[ aliensh 2.1.18a (C) ARDA/Alice: Andreas.Joachim.Peters@cern.ch/Derek.Feichtinger@cern.ch]
aliensh:[alice] [1] /alice/cern.ch/user/r/rvernet/ >ls -l Analysis/
drwxr-xr-x  rvernet  rvernet           0 Jan 17 17:14  .runAnalysis.jdl
-rwxr-xr-x  rvernet  rvernet        387 Jan 17 16:29  CreateChainFromXML.C
-rwxr-xr-x  rvernet  rvernet        396 Jan 17 23:51  mergerootfile-sequential.jdl
drwxr-xr-x  rvernet  rvernet           0 Jan 17 17:47  output
-rwxr-xr-x  rvernet  rvernet        632 Jan 17 17:14  runAnalysis.jdl
-rwxr-xr-x  rvernet  rvernet       1746 Jan 17 13:52  sample.xml
-rwxr-xr-x  rvernet  rvernet       1797 Jan 17 16:12  strangeAna.C
aliensh:[alice] [2] /alice/cern.ch/user/r/rvernet/ >submit Test/cream_aliroot.jdl
Submit submit Test/cream_aliroot.jdl
submit: Your new job ID is 41931762
aliensh:[alice] [3] /alice/cern.ch/user/r/rvernet/ >ps
 rvernet 41552696  -- D      aliroot
 rvernet 41552697  -- D      date
 rvernet 41552698  -- D      aliroot
 rvernet 41552699  -- D      date
 rvernet 41552927  -- D      aliroot
 rvernet 41552928  -- D      date
 rvernet 41552929  -- D      aliroot
 rvernet 41552940  -- D      date
 rvernet 41891047  -- D      aliroot
 rvernet 41894985  -- D      aliroot
 rvernet 41931761  0 I      aliroot
 rvernet 41931762  0 I      aliroot
aliensh:[alice] [4] /alice/cern.ch/user/r/rvernet/ >whereis /alice/data/2009/LHC09d/000104160/ESDs/pass2/09000104160018.10/root_archive.zip
Feb 11 12:54:00 info The file LHC09d/000104160/ESDs/pass2/09000104160018.10/root_archive.zip is in
SE => ALICE::CCIN2P3::SE pfn =>root://ccxrdsn038.in2p3.fr:1094//07/23620/8e0b747a-f759-11de-a909-0018fe730ae5

SE => ALICE::FZK::SE pfn =>root://f01-120-123-e.gridka.de:1094//07/23620/8e0b747a-f759-11de-a909-0018fe730ae5

SE => ALICE::Subatech::SE pfn =>root://nanxrdmgr01.in2p3.fr:1094//07/23620/8e0b747a-f759-11de-a909-0018fe730ae5

SE => ALICE::JINR::SE pfn =>root://lcgxrdr01.jinr.ru:1094//07/23620/8e0b747a-f759-11de-a909-0018fe730ae5
aliensh:[alice] [5] /alice/cern.ch/user/r/rvernet/ >
```





# Job submission



- Jobs submitted via AliEn shell
  - ◆ Job optimizer selects the best sites that satisfy the request
  - ◆ Performs splitting into subjobs
  - ◆ The subjobs will run where the data is and only there !! *(in principle)*
    - NB: several sites may contain the same data (mirror)
- Job Agents are launched on different sites
  - ◆ Check availability of site
  - ◆ Check the required software
    - If not, the package manager will install it
  - ◆ Pulls 'real' jobs and runs them on the WN
- Which VObox (group of CEs) to use?
  - ◆ If not specified, the Job Optimizer decides itself
  - ◆ This can be specified in the JDL
    - `Requirements=other.CE=="Alice::CCIN2P3::CCIN2P3-CREAM"`



# What happens at CCIN2P3?



- T1+T2 site → we do many different things
  - ◆ Raw data backup (HPSS)
  - ◆ Reconstruction (jobs 'aligrid')
  - ◆ MC production (jobs 'aligrid')
  - ◆ Scheduled analysis (from production, jobs 'aligrid')
    - Creation of "Physics" data → ESD, AOD
  - ◆ Chaotic (end-user) analysis
    - Analysis of ESD,AOD (jobs 'alicexxx')
- 2 computing elements
  - ◆ LCG CE → VObox `cc1cgalice01` submits on it
  - ◆ CREAM CE → VObox `cc1cgalice02` submits on it
    - More promising, will be the default at some point
    - `cc1cgalice01` will become a backup VObox
- Storage
  - ◆ Tape (HPSS) : raw data transfers (FTS) from T0
  - ◆ Disk (xrootd) : output from reconstruction done in T0 (→ send via FTS) and done in T1 (CC)
- Data staging Tape → Disk
  - ◆ Feature not yet implemented
  - ◆ Wish to use TReqs !





# Job behaviour @ CC



- Common software for all jobs in
  - /afs/in2p3.fr/grid/toolkit/alice
  - 1 AFS volume → stress ?
  - Resource used with BQS: `u_stress_afs_alice`
    - Limit the incoming job flow
- Memory issue on production jobs
  - Simulation & reconstruction require large amount of memory
    - Min. 2GB
  - We reached a value of 4.5GB for some jobs
    - Need to use the BQS J class for these jobs (hopefully temporary issue)
  - AliRoot offline core goal is to make it get down to 2GB again
  - However we should take an action soon



# Storage @ CC



- RAW data replication
  - ◆ Should do CASTOR@CERN → HPSS@CCIN2P3 via FTS
  - ◆ Staging HPSS → xrootd not in place yet
  - ◆ ⇒ no raw data at CCIN2P3 at the moment
- Subsequent reconstruction passes have been stored on xrootd@CCIN2P3
  - ◆ → user analysis can be performed on ESD/AOD
- Simulation
  - ◆ Performed partly at CCIN2P3
  - ◆ → everything on xrootd
- SE Monitoring
  - ◆ Files are regularly sent and stored in SE's
  - ◆ Provides “working status” of storage services





# Storage issues



- Deployment procedure of xrootd software with ALICE authentication library
- Effective data transfer rate from xrootd servers
  - “GRID SE” slower than “LAF” server
  - Is this expected ?
  - Software or bandwidth issue ?



# PROOF cluster (LAF)



- Computing model of ALICE and Analysis Framework designed to work on PROOF
  - ◆ Exact same code must run on local, GRID and PROOF
- See LAF tutorial @ CC (last week) for more info





# Current status and prospects

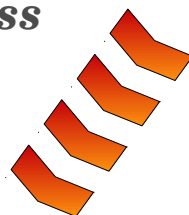


ccli (client)

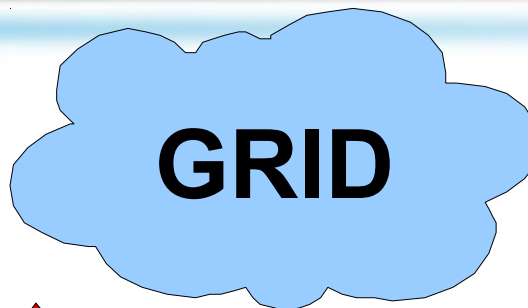


query

*Direct access  
to our SE !*



**GRID**



*Manual staging*



*User friendly  
data staging  
and  
management*

**PROOF**



1 master



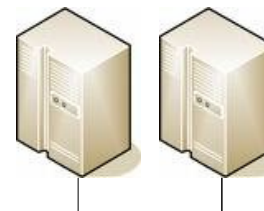
20 slaves \* 8 cores



Scientific Linux



**Data !  
(xrootd  
protocol)**



xrootd server



32TB



questions?



# FTS vers SE xrootd (december 17<sup>th</sup>-18<sup>th</sup>)

SEs average transfer rates

