

ALICE @ ARDA

P. Cerello – INFN Torino

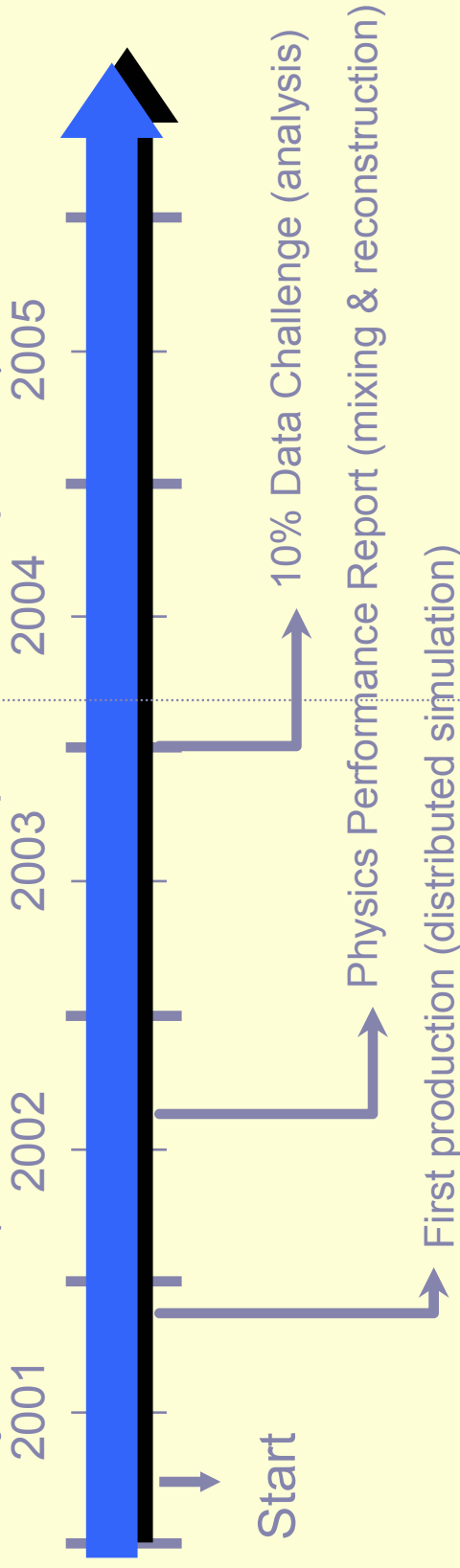
ARDA Workshop

June 23 2004

The ALICE Approach (ALiEn)

- There are millions lines of code in OS dealing with GRID issues
- Why not using them to build the minimal GRID that does the job?
 - Fast development of a prototype, can restart from scratch etc etc
 - Hundreds of users and developers
 - Immediate adoption of emerging standards

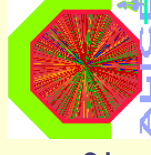
• ALiEn by ALICE (5% of code developed, 95% imported)



Functionality + Interoperability + Performance, Scalability, Standards
+
Simulation + Reconstruction + Analysis

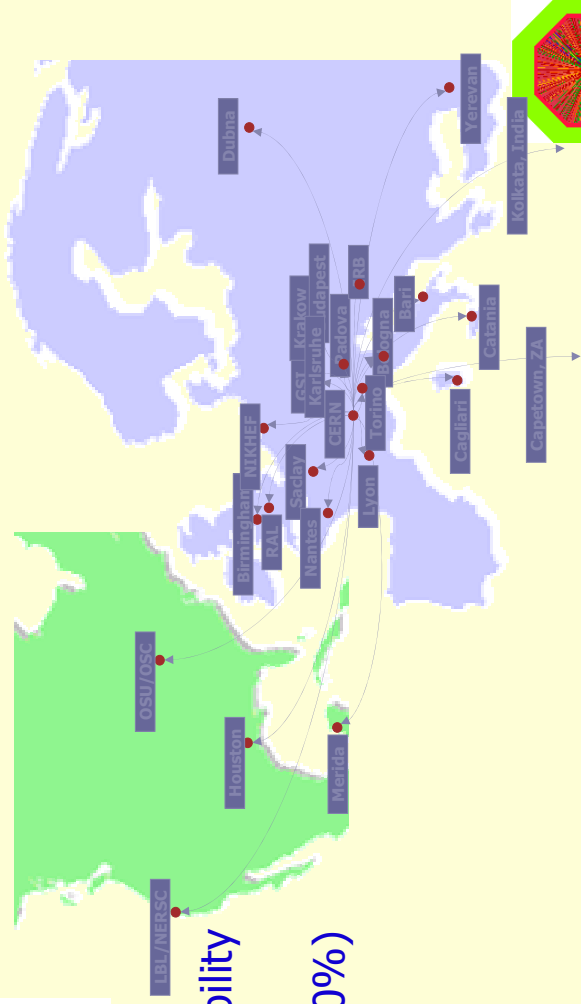
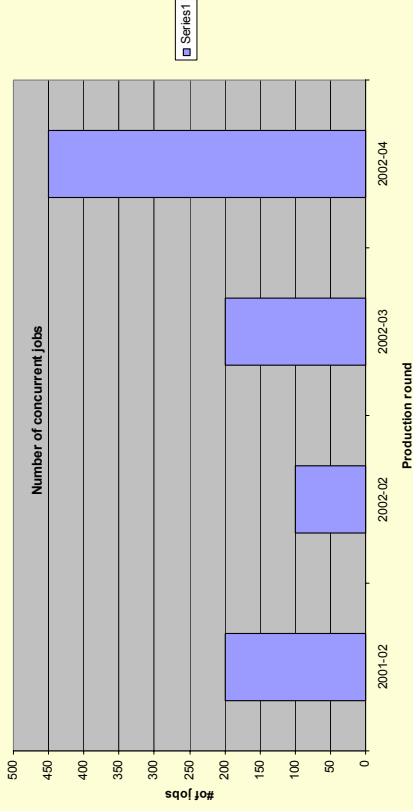
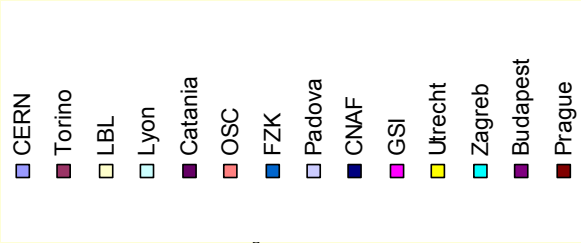
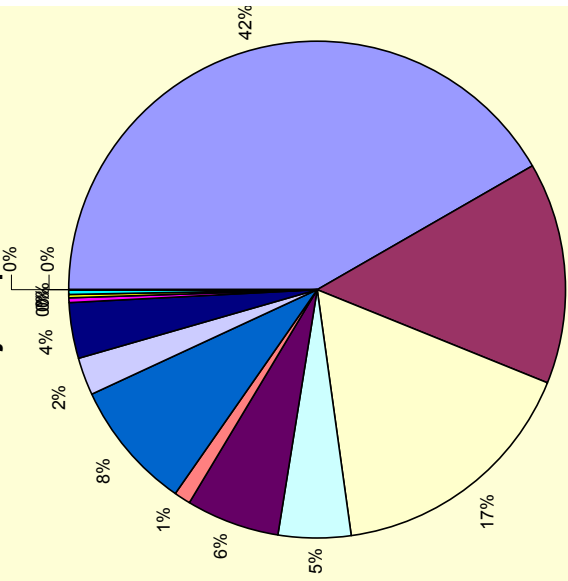
June, 23rd, 2004

ARDA Workshop

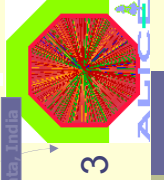


AliEn activity for the PPR

Total jobs per site



- ◆ 32 sites configured
- ◆ 5 sites providing mass storage capability
- ◆ 12 production rounds
- ◆ 22773 jobs validated, 2428 failed (10%)
- ◆ Up to 450 concurrent jobs
- ◆ 0.5 operators

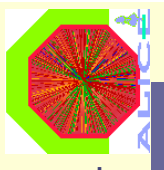
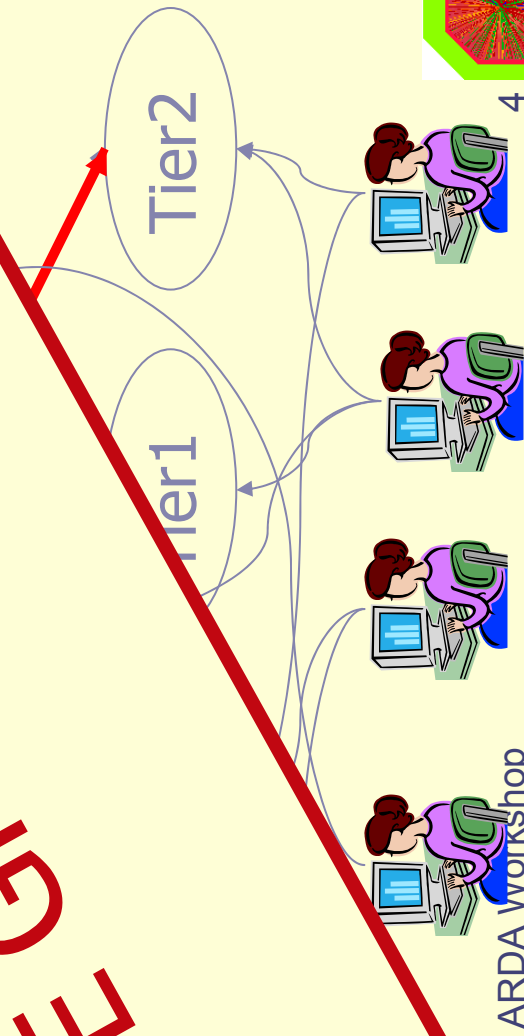


PDC 3 sch

En job control
a transfer

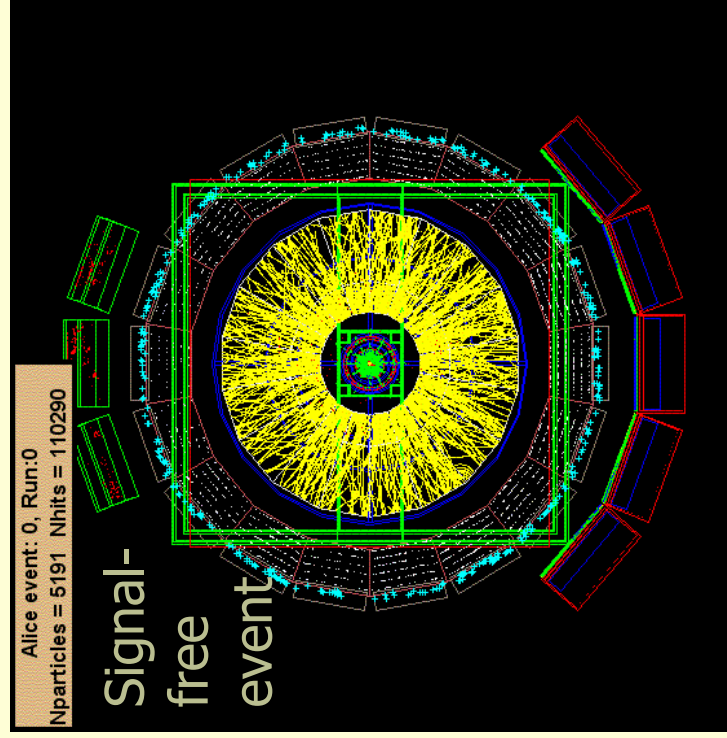
Production of RAW
Shipment of RAW +
Reconstruction
Analysis

NO MORE
DO IT ALL
THE GRIND

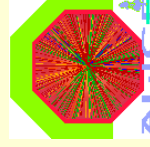
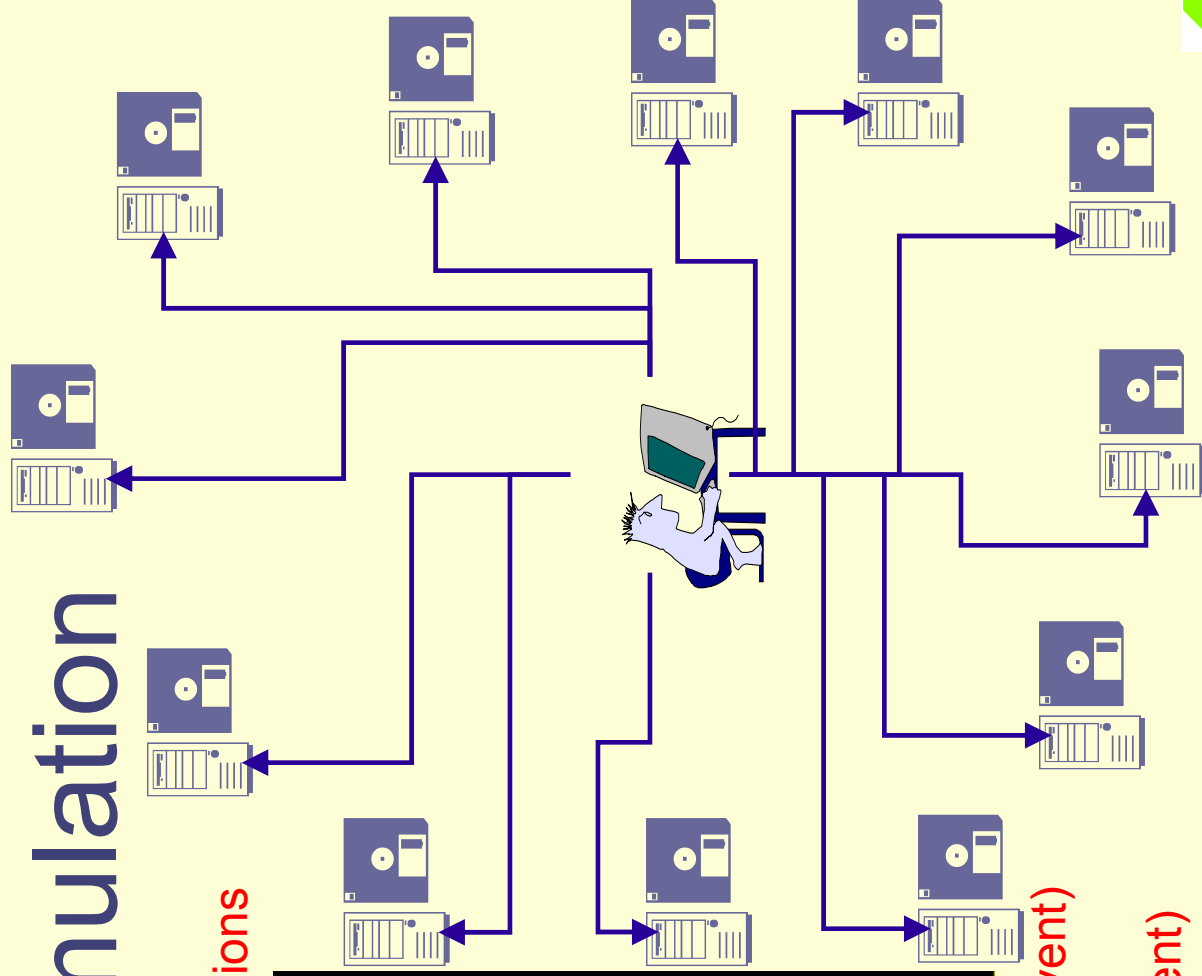


Simulation

Small input with interaction conditions

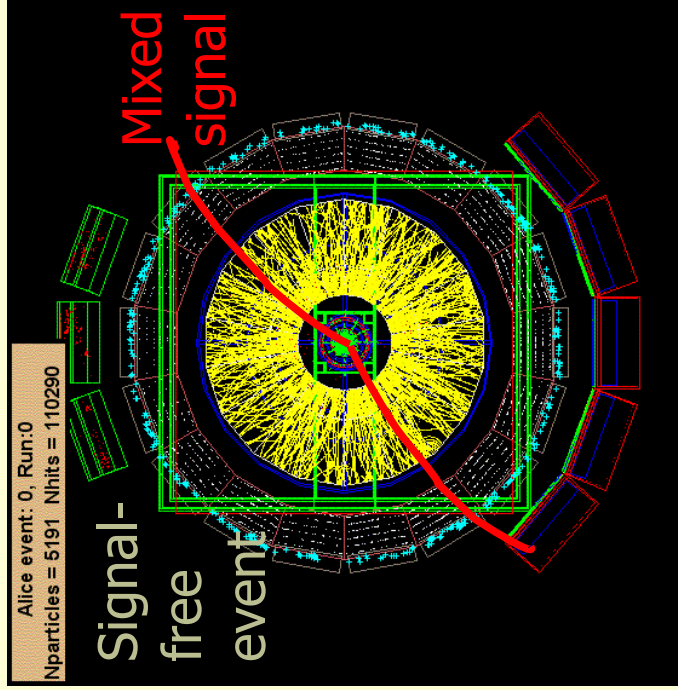


Large distributed output (1 GB/event)
with simulated detector response
Long execution time (10 hours/event)

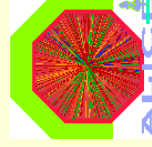
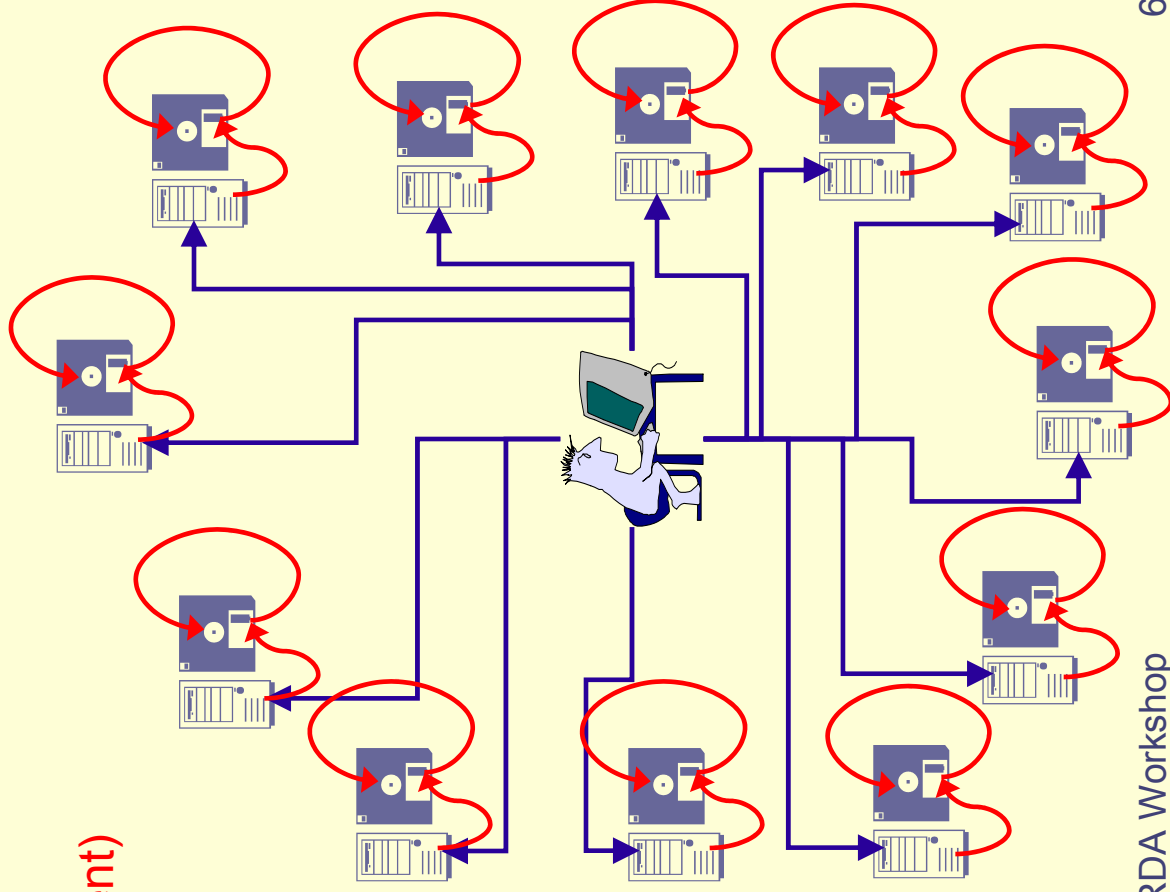


Merging + Reconstruction

Large **distributed** input (1 GB/event)

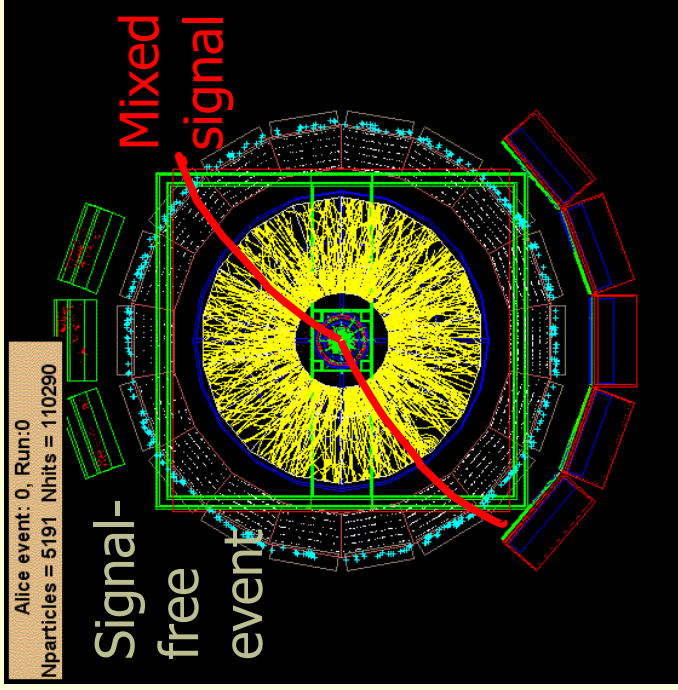


Fairly Large **distributed** output (100 MB/event, 7MB files) with reconstructed events

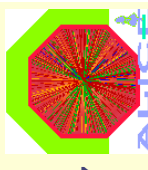
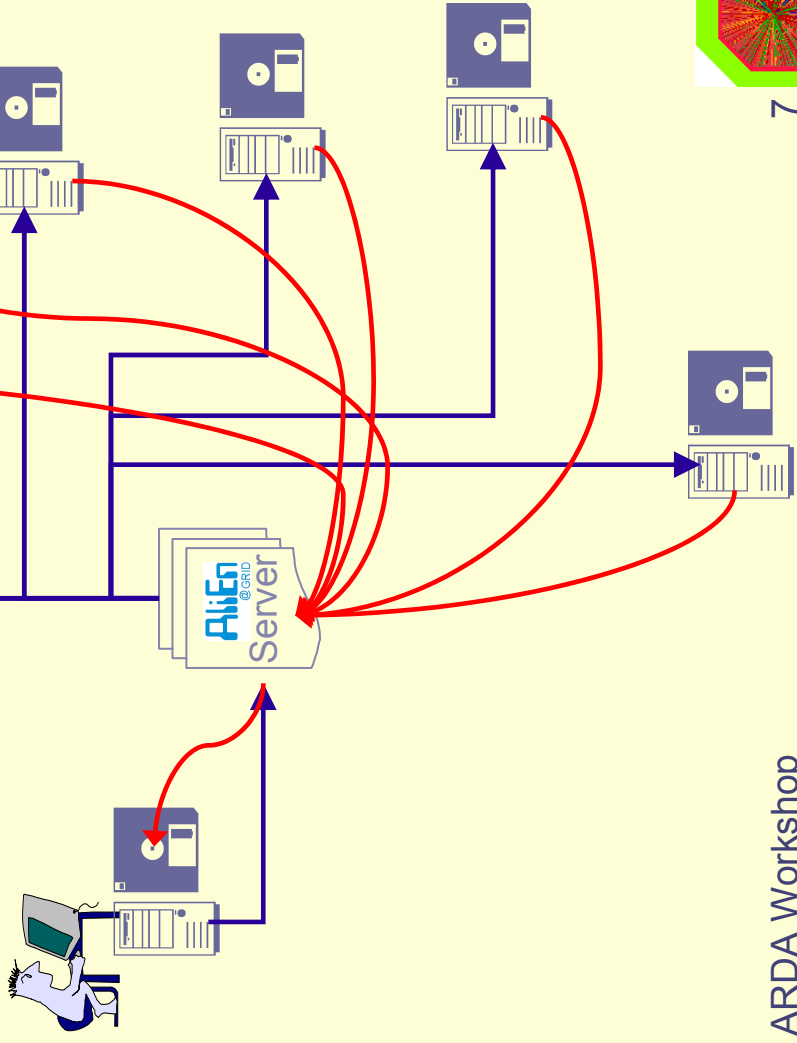


Analysis

Large distributed input (100 MB/event)



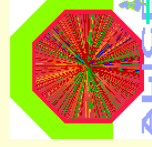
Fairly small merged output



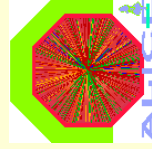
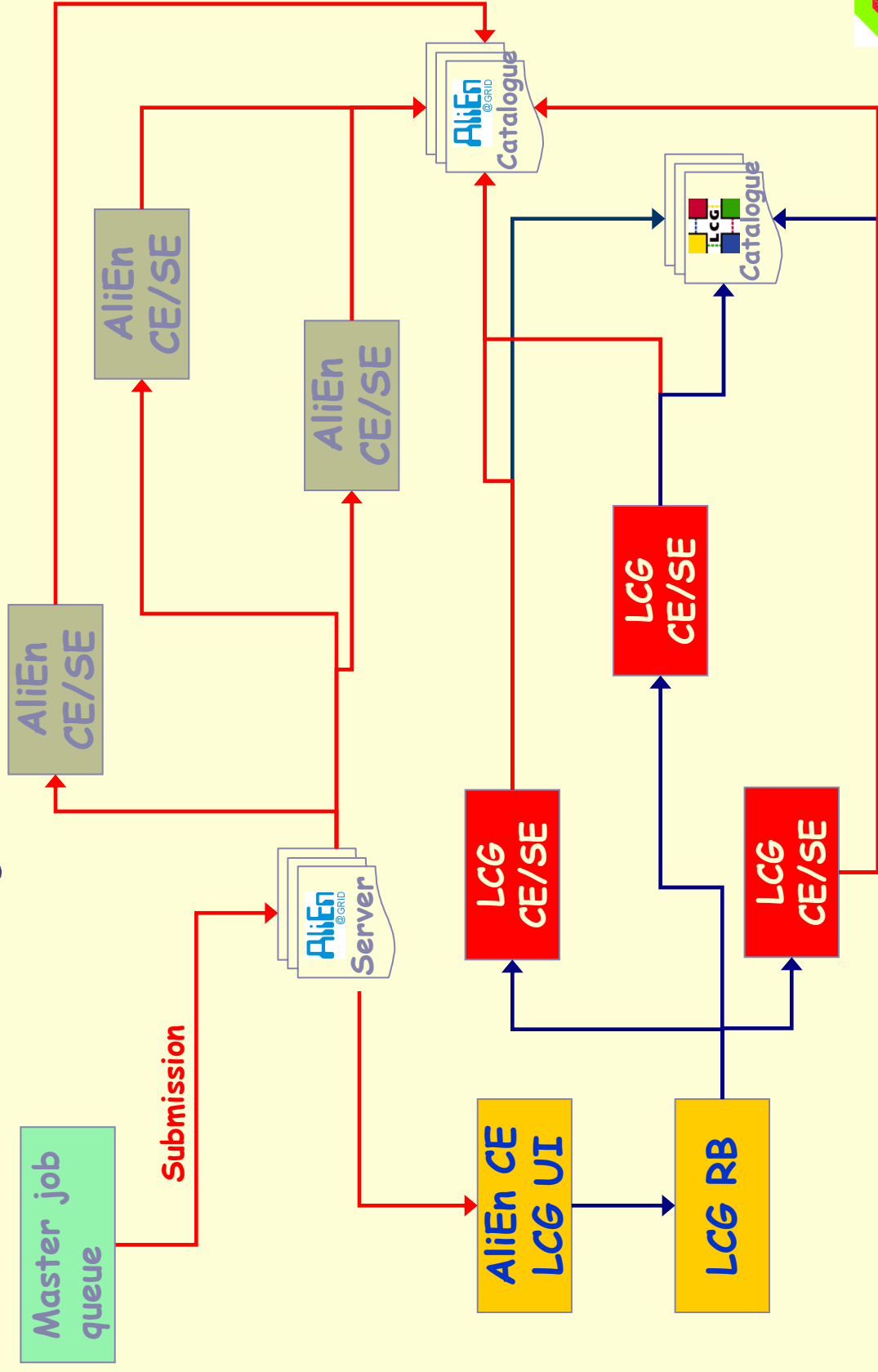
Phases of ALICE Physics Data Challenge

2004

- Phase 1 - production of underlying events using heavy ion MC generators
 - Status – 100% complete in about 2 months of active running
 - Basic statistics - ~ 1.3 million files, 26 TB data volume
- Phase 2 – mixing of signal events in the underlying events
 - Status – about to begin
- Phase 3 – analysis of signal+underlying events:
 - Goal – to test the data analysis model of ALICE
 - Status – will begin in 2 months

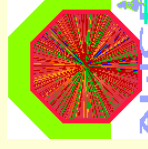


DC2004 layout: a "Meta-Grid"



Design strategy

- “Pull-model” well suited for implementing high-level submission systems, since it does not require knowledge about the periphery, that may be very complex
- **Use AliEn as a single general front-end**
 - Owned and shared resources are exploited transparently
- **Minimize points of contact and don't be invasive**
 - No need to reimplement services etc.
 - No special services required to run on remote CE/WNs
- **Make full use of provided services: Data Catalogues, scheduling, monitoring...**
 - Let the Grids do their jobs (they should know how)
- **Use high-level tools and APIs to access Grid resources**
 - Developers put a lot of abstraction effort into hiding the complexity and shielding the user from implementation changes

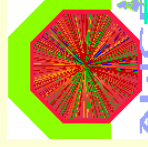


Implementation

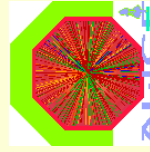
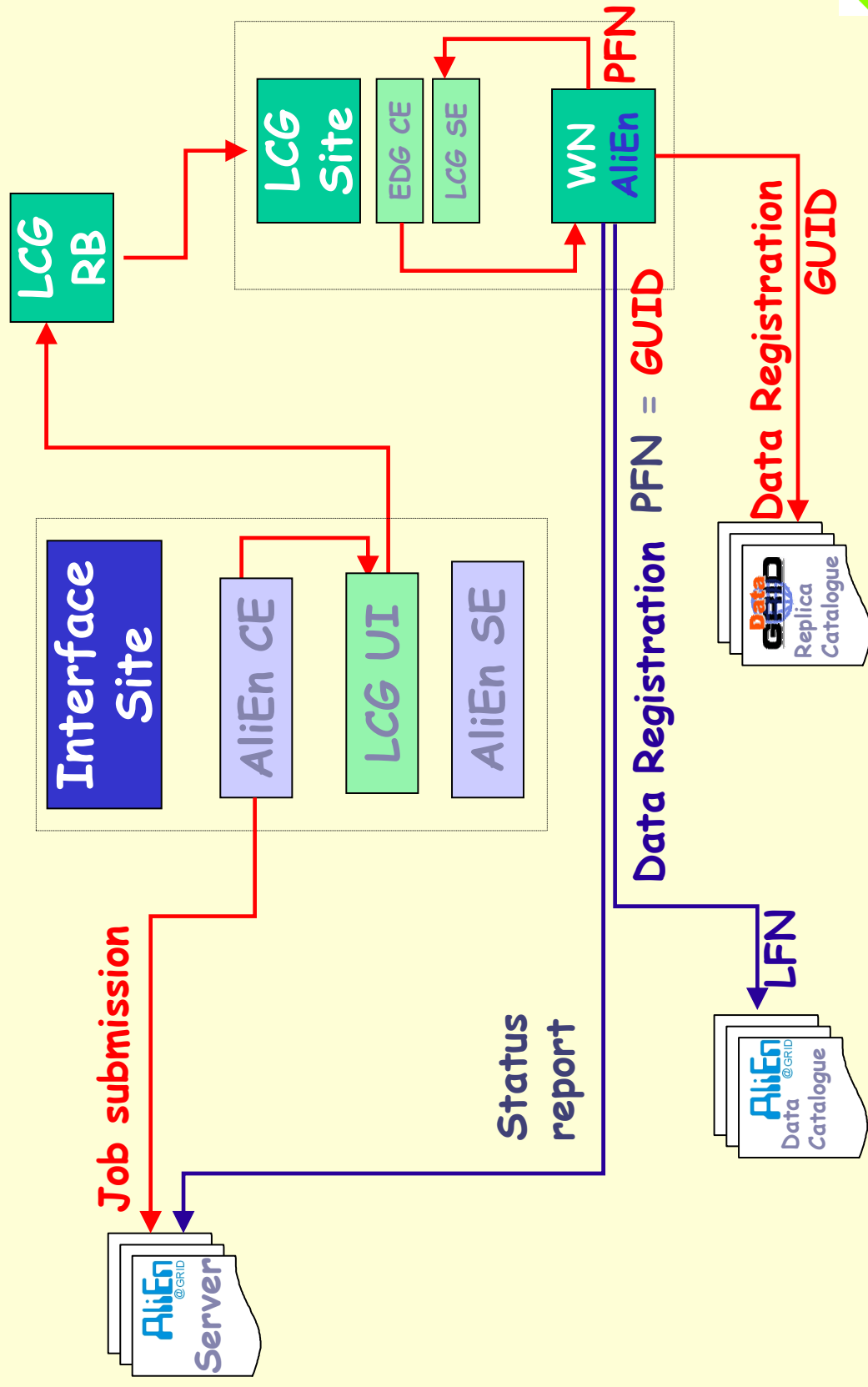
The whole of LCG computing is seen as a single, large AliEn CE associated with a single, large SE

Implementation: manage LCG (& Grid.it) resources through a “gateway”: an AliEn client (CE+SE) sitting on top of an LCG User Interface

- **Job management interface**
 - JDL translation (incl. InputData statements)
 - JobID bookkeeping
 - Command proxying (Submit, Cancel, Status query...)
- **Data management interface**
 - Put() and Get() implementation
 - AliEn PFN is LCG GUID
 - (plus SE host to allow for optimisation)
 - AliEn knows nothing about LCG replicas (but RLS/ROS do!)

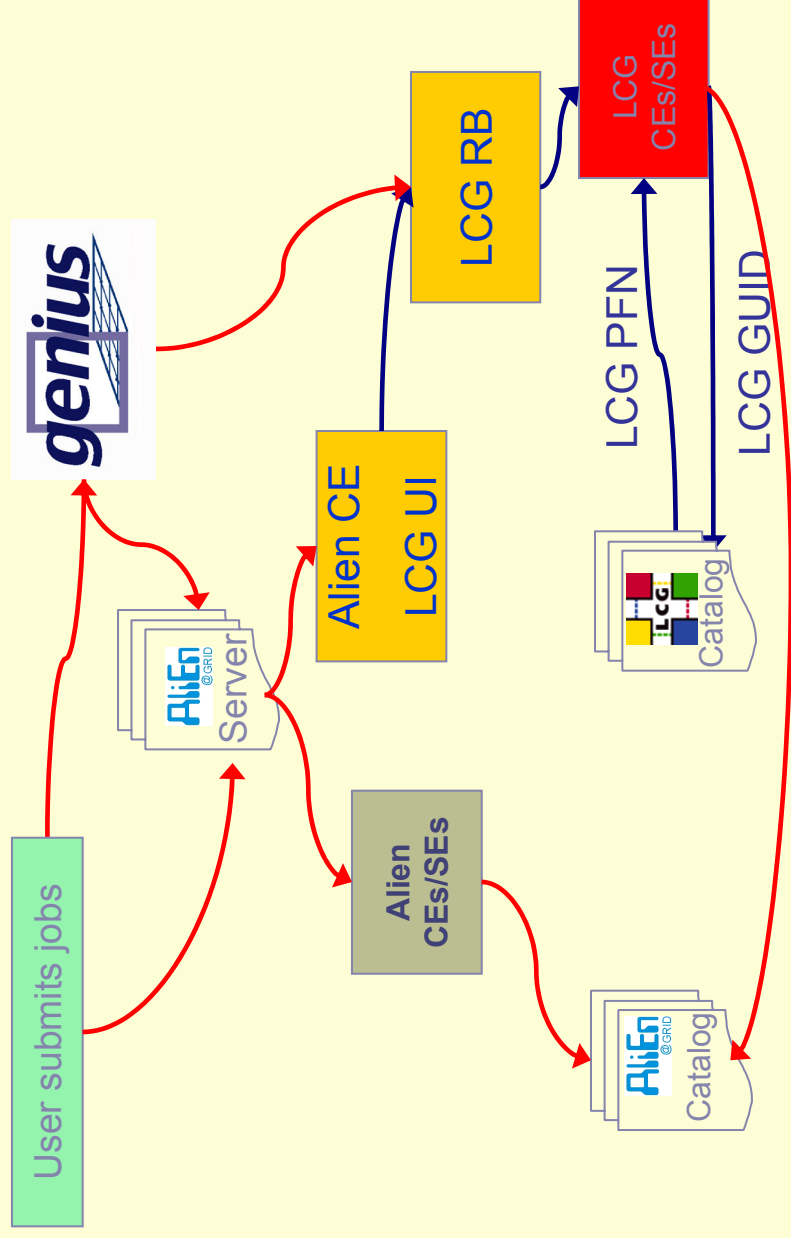


Interfacing AliEn and LCG



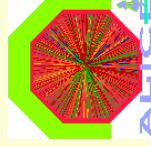
AliEn, Genius & EDG/LCG

- LCG-2 is one CE of AliEn, which integrates LCG and non LCG resources
 - If LCG-2 can run a large number of jobs, it will be used heavily
 - If LCG-2 cannot do that, AliEn selects other resources, and it will be less used

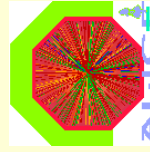
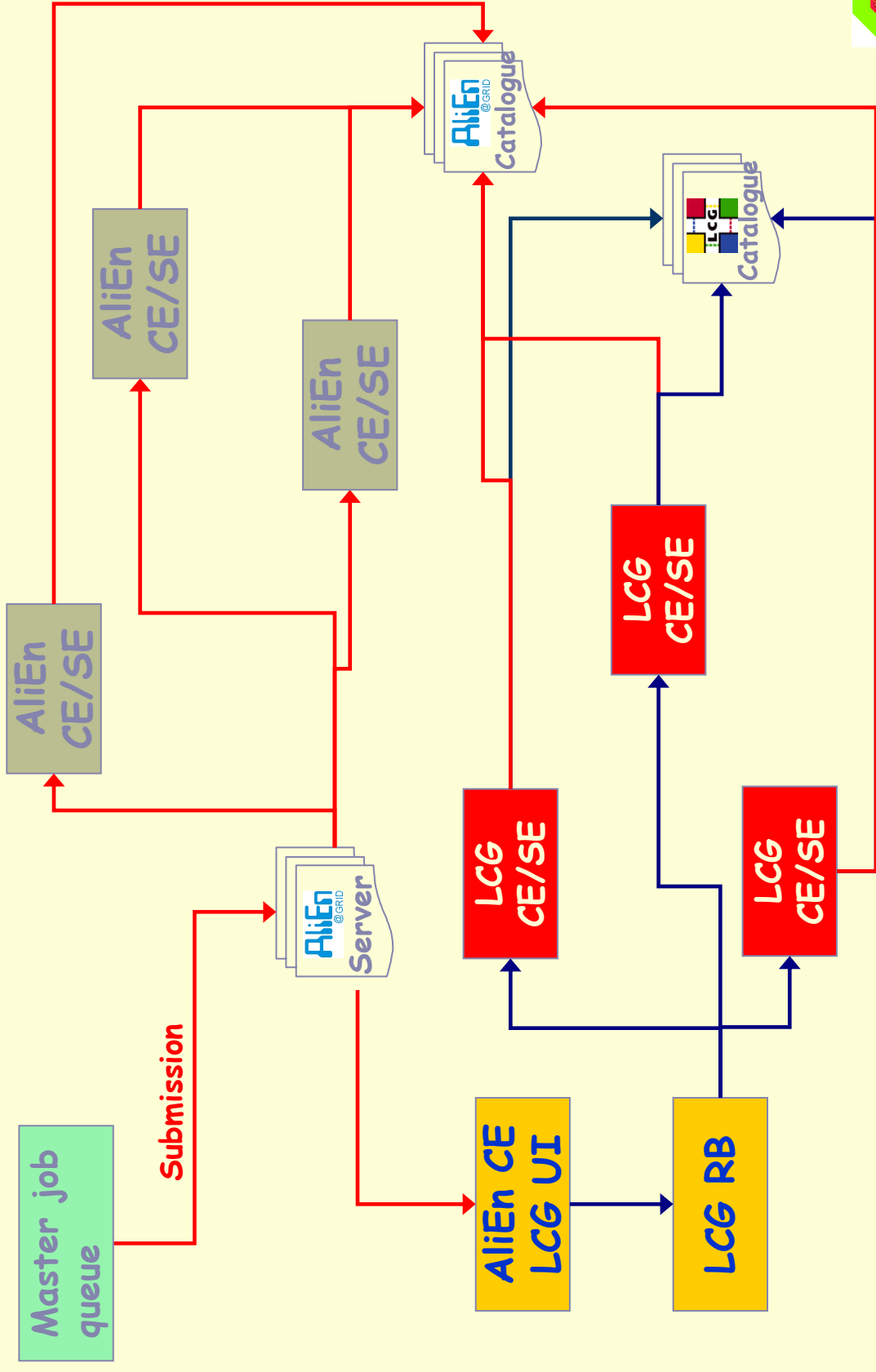


LCG GUID = AliEn PFN
ARDA Workshop

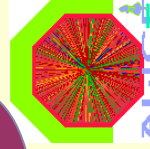
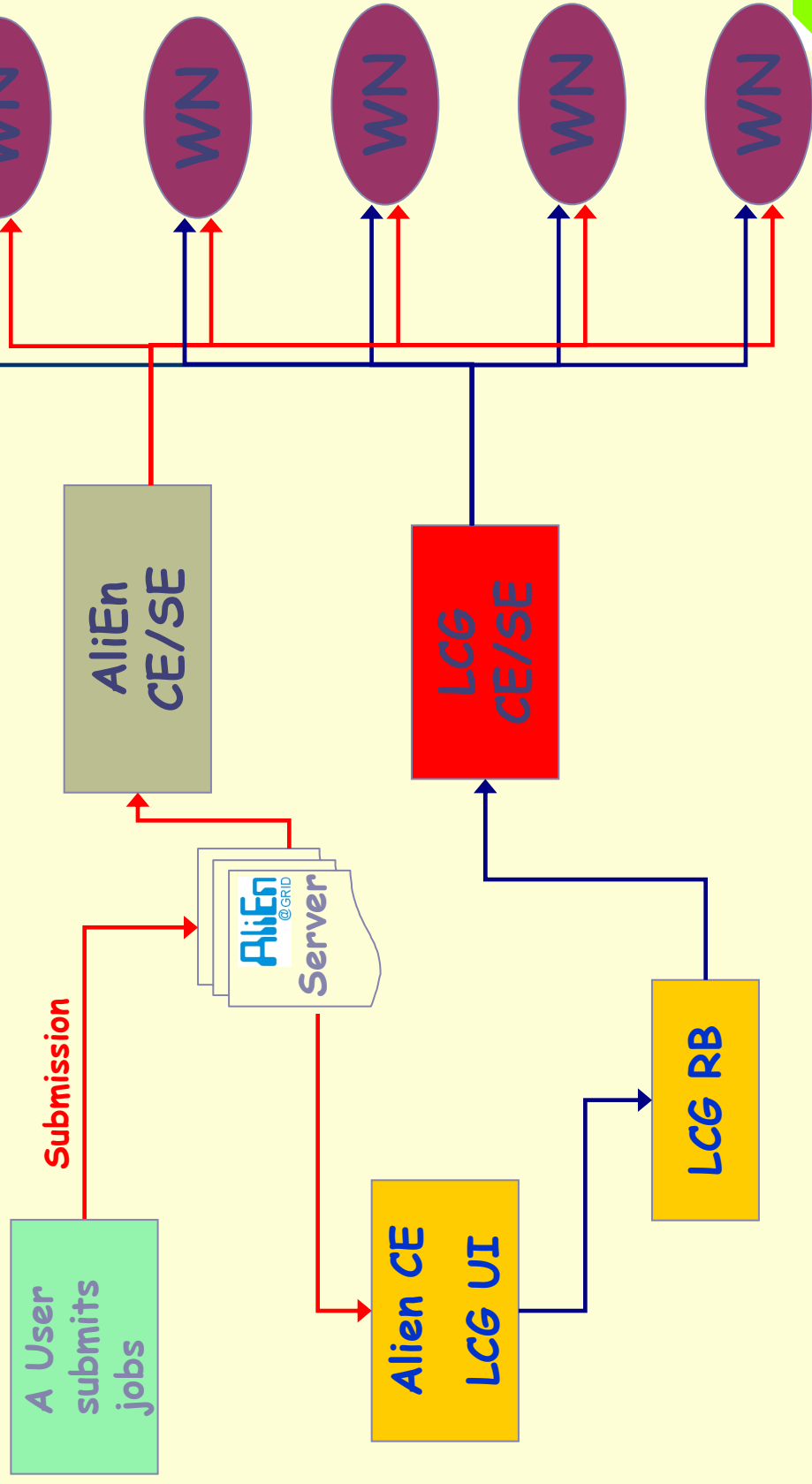
June, 23rd, 2004



DC2004 layout: a "Meta-Grid"



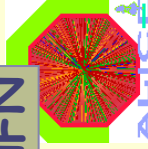
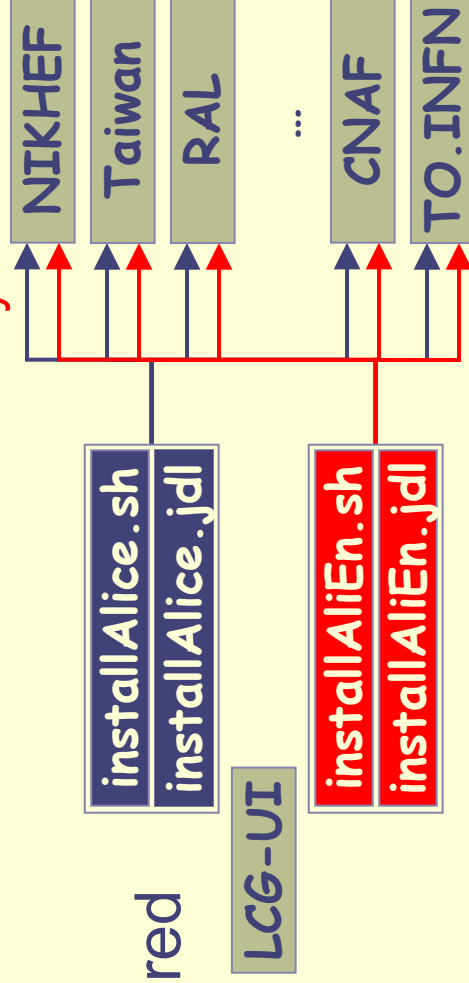
Double Access (CNAF & ct.infn.it)



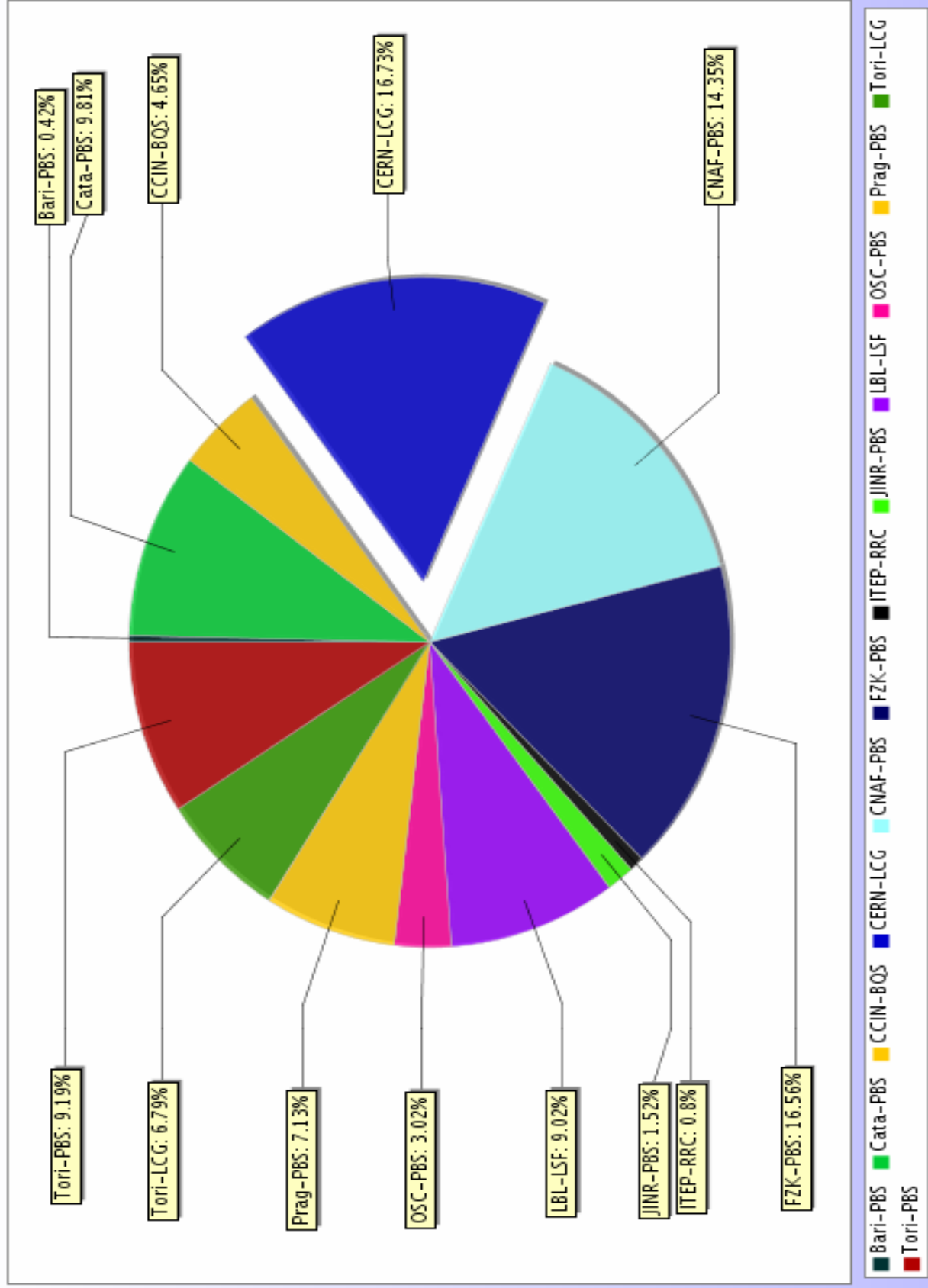
Software installation on LCG

- Both AliEn and AliRoot installed via LCG jobs
 - Do some checks, download tarballs, uncompress, build environment script and publish relevant tags
 - Single command available to get the list of available sites, send the jobs everywhere and wait for completion. Full update on LCG-2 + GRID.IT (16 sites) takes ~30'
 - Manual intervention still needed in few sites (e.g. CERN/LSF)
 - Ready for integration into AliEn automatic installation system

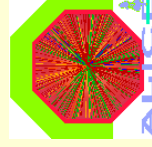
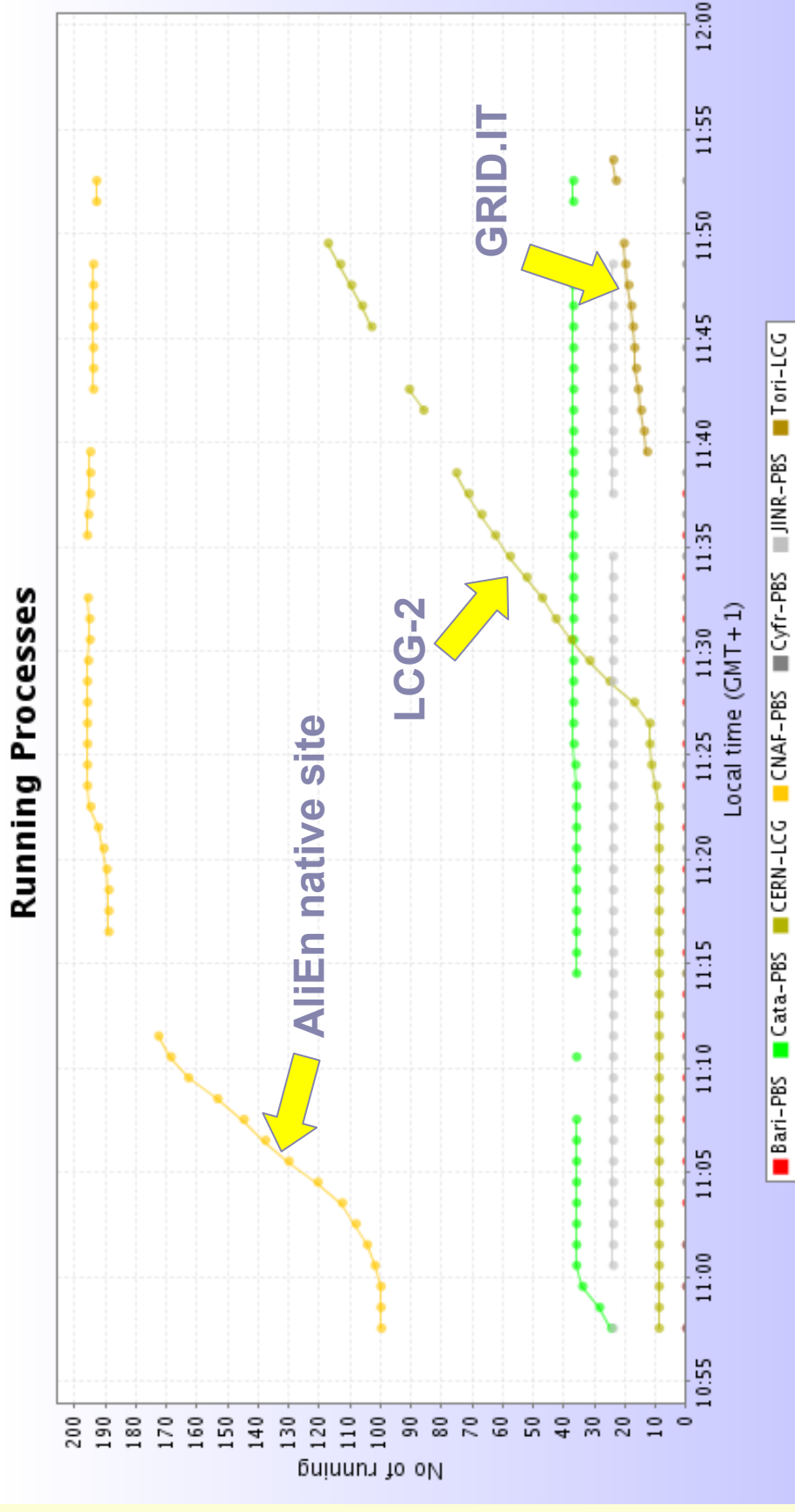
- Experiment software shared area misconfiguration caused most of the trouble in the beginning



Jobs done

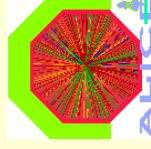
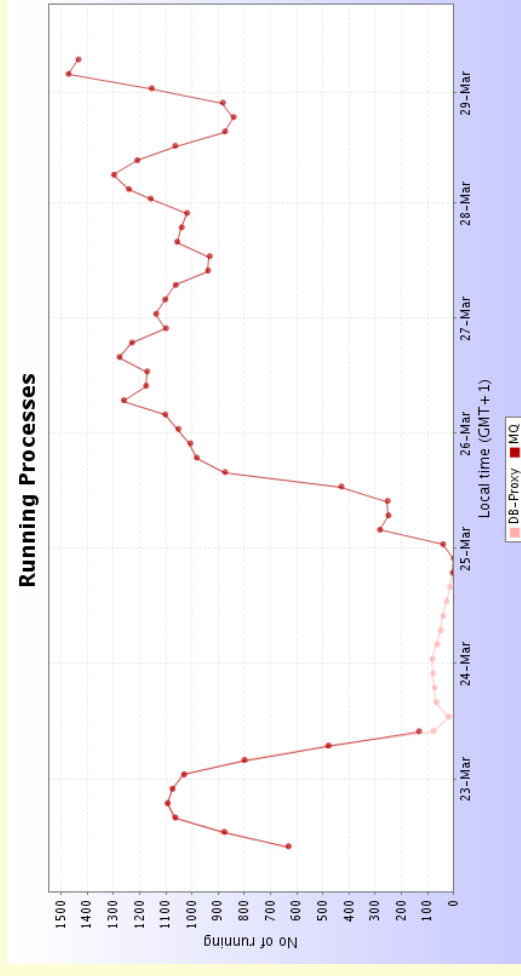
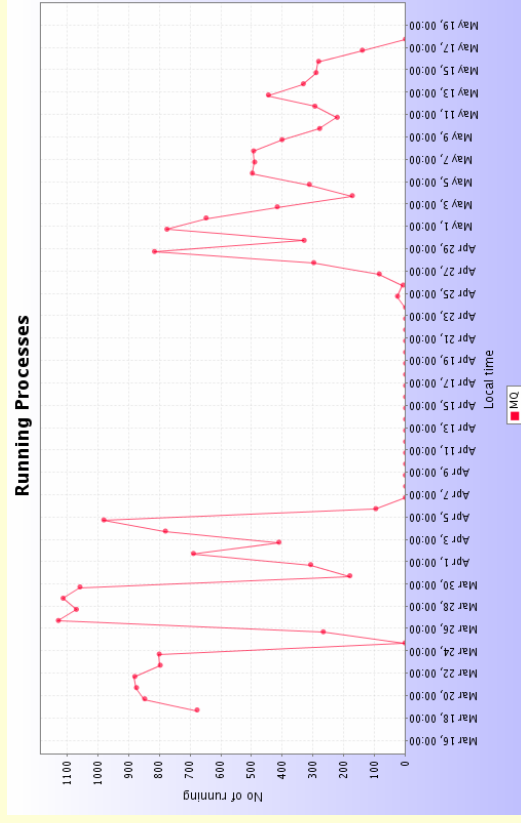


DC monitoring is helpful



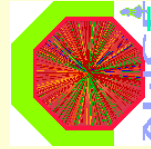
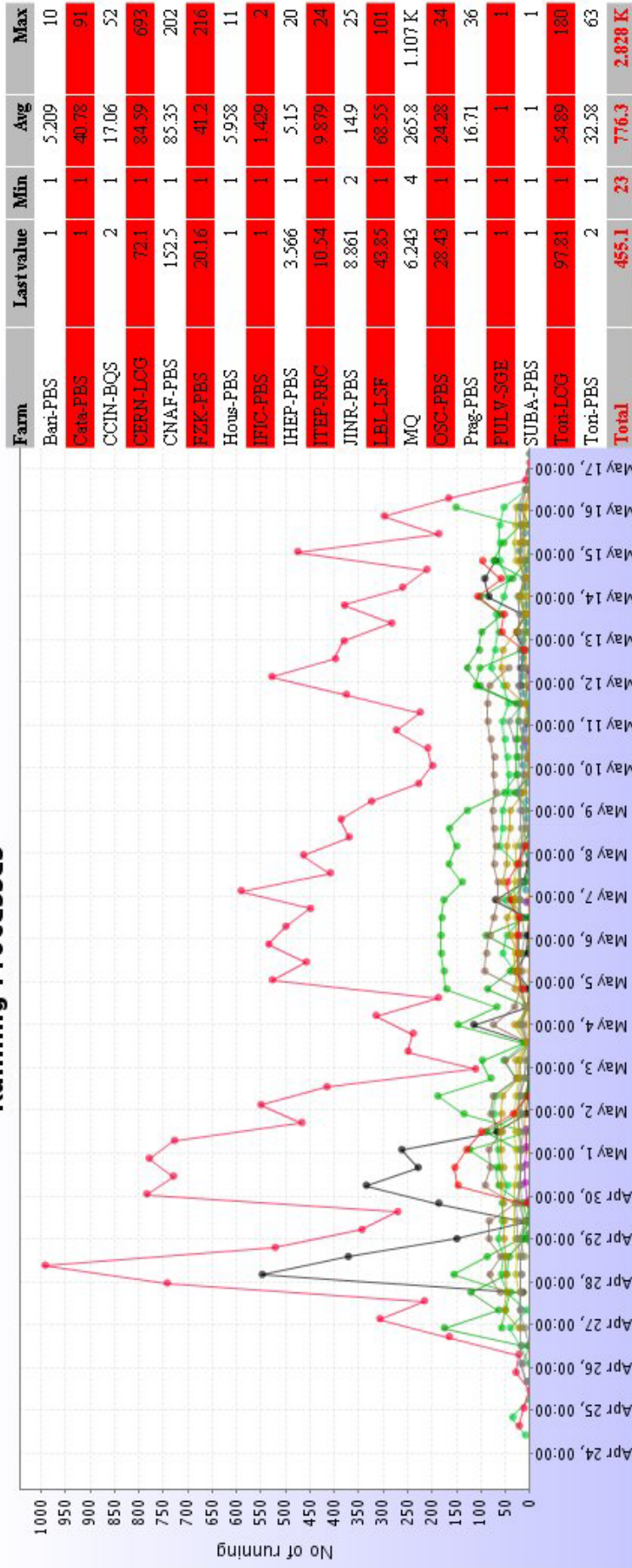
DC Monitoring

- Total CPU profile over the first phase of the DC:
 - Aiming for continuous running, not always possible due to resources constraints



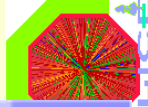
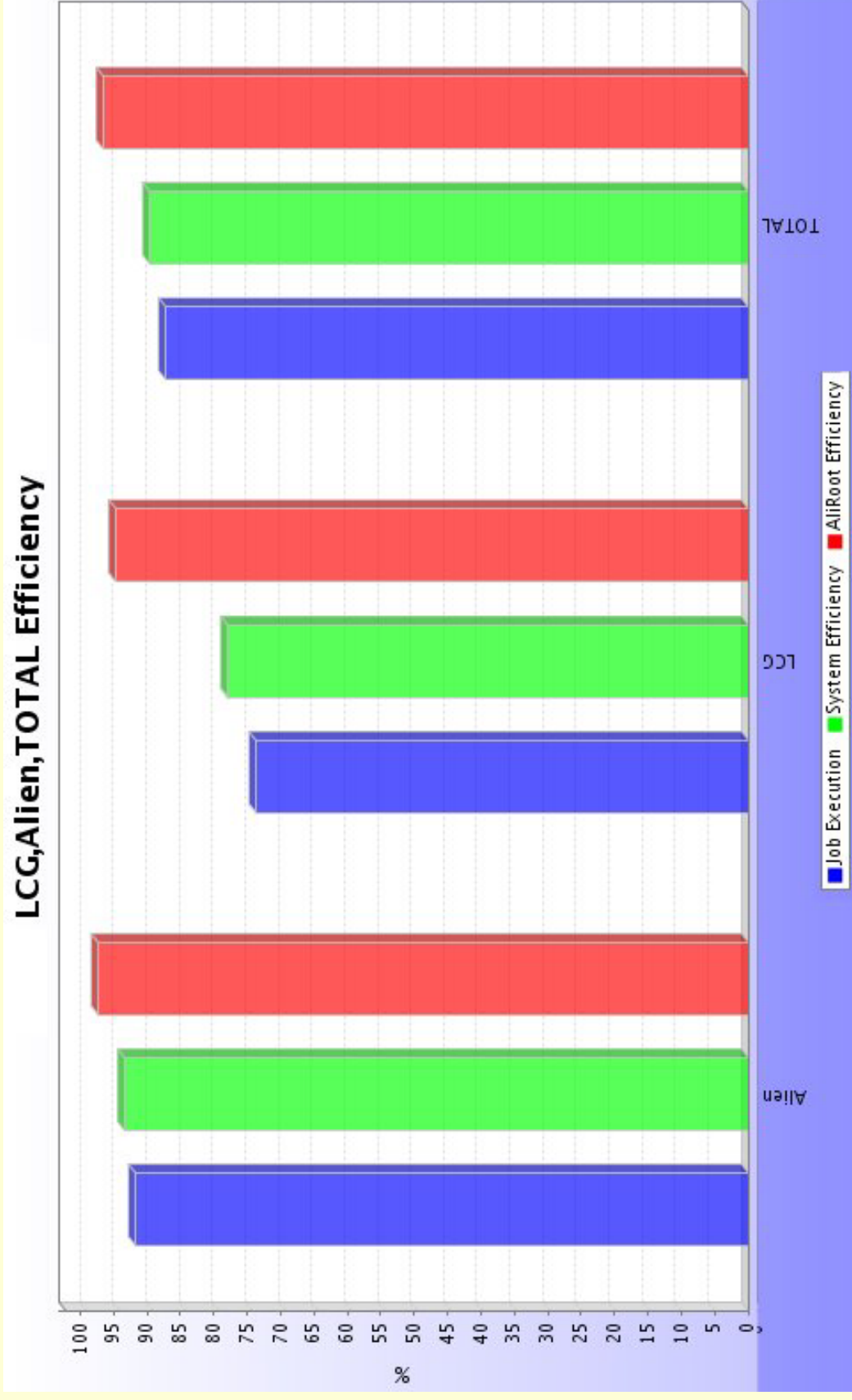
Running Job History

Running Processes



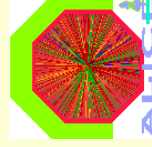
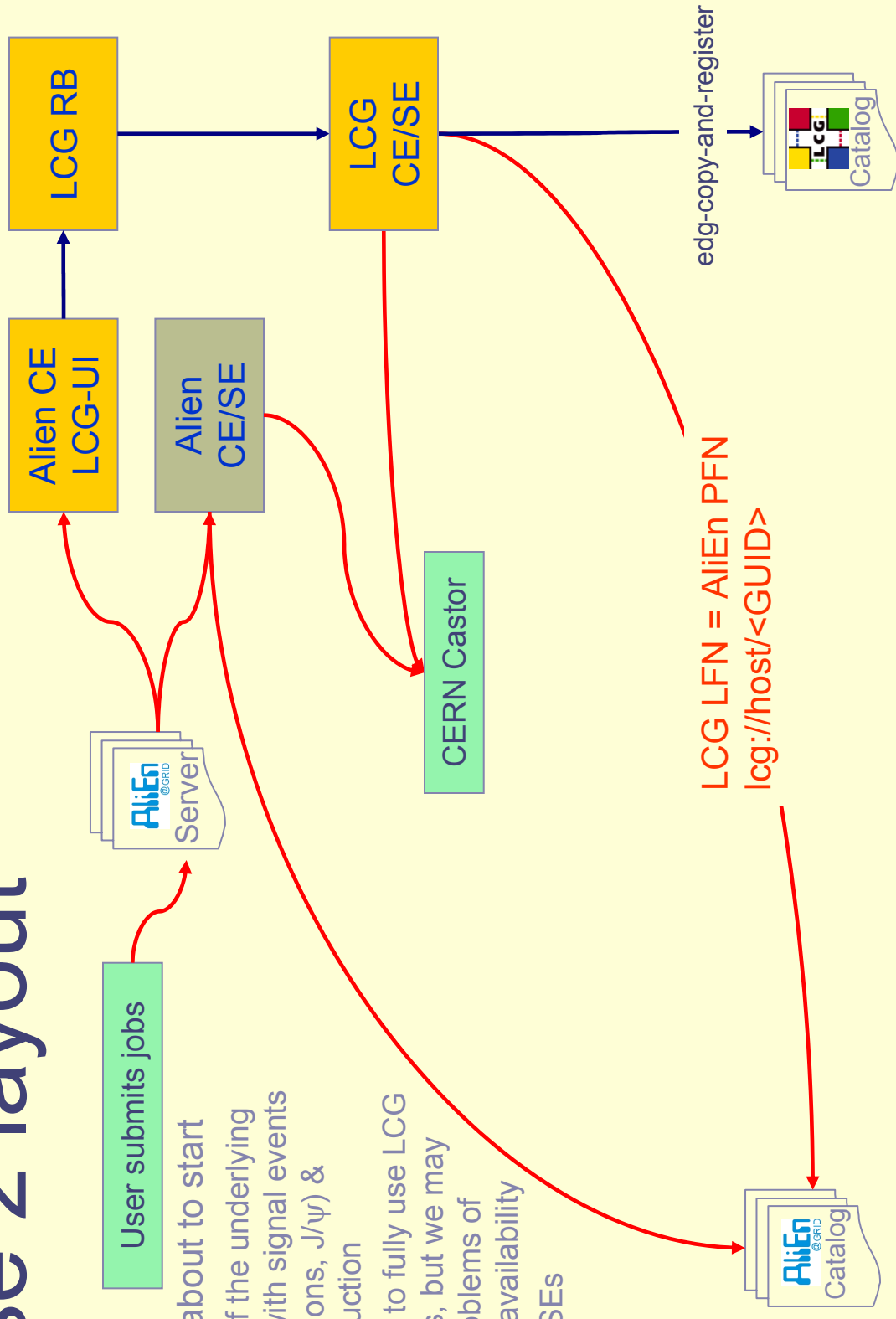
Groups Efficiency

- Job execution, System and AliRoot efficiency per groups



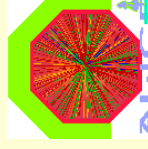
Phase 2 layout

- Phase 2 -- about to start
 - Mixing of the underlying events with signal events (jets, muons, J/ψ) & reconstruction
 - We plan to fully use LCG DM tools, but we may have problems of storage availability at local SEs

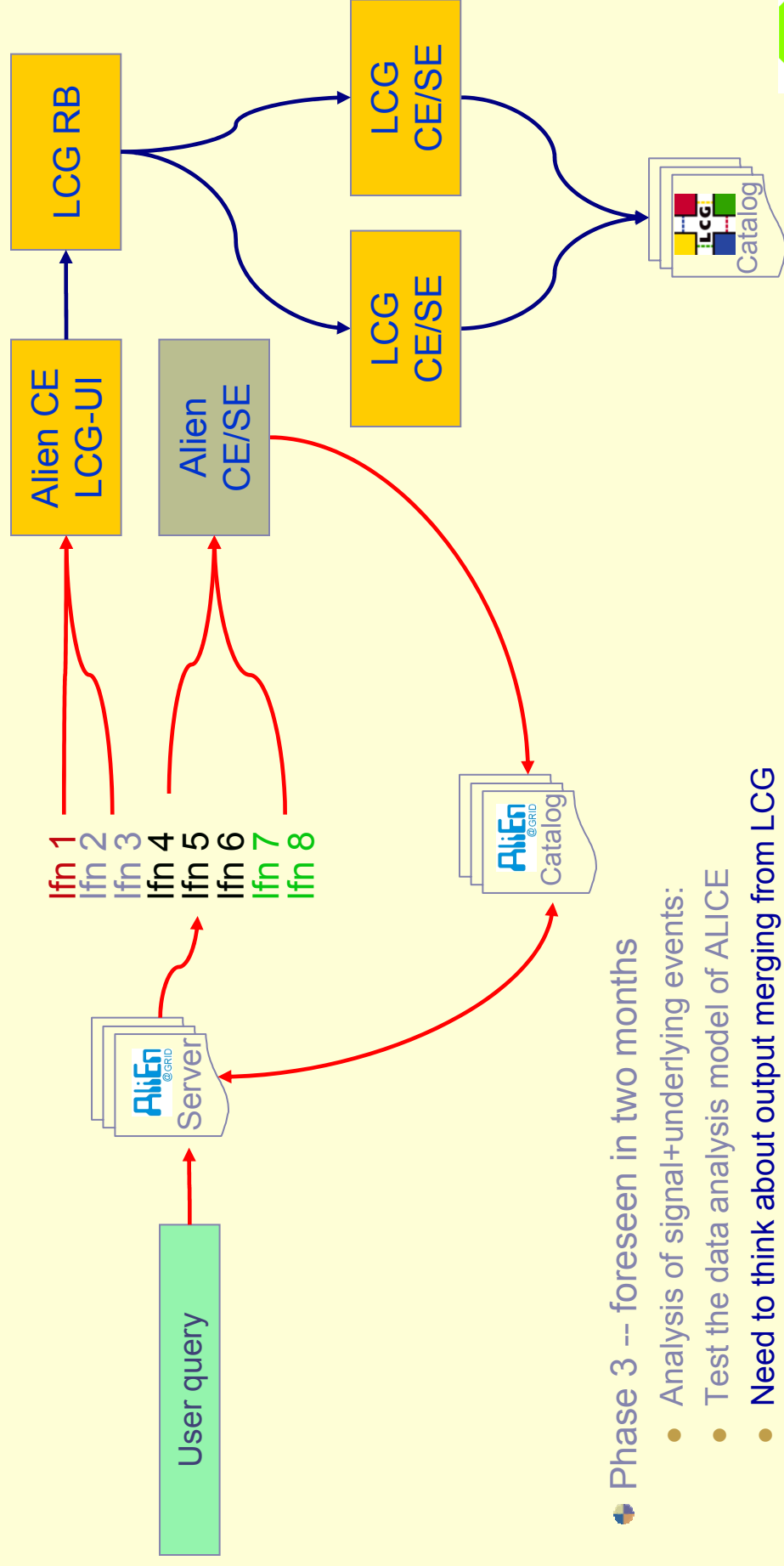


Issues with Phase 2

- Phase II will generate lots (1M) of (rather small ~7MB) files
- It will start in a few days
- We are testing a plug-in to AliEn using tar to bunch small files
- The space available on the some LCG SEs seems very small... we might be able to run just for few hours before shutting them off
- Preparation of the LCG-2 JDL is more complicated, due to the use of the data management features



Phase 3 layout

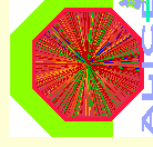


Phase 3 -- foreseen in two months

- Analysis of signal+underlying events:
- Test the data analysis model of ALICE
- **Need to think about output merging from LCG**
- **ARDA availability would simplify phase 3!!!!**

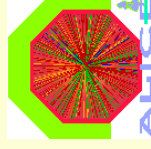
June, 23rd, 2004

ARDA Workshop



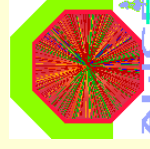
Considerations I

- AliEn command-line tools OK for DC running and resources control
- AliEn quick feedback from the CE and WN proved to be essential for early spotting of problems
 - Jobs are constantly monitored through heartbeat
 - Failure thresholds are put in place to close problematic CEs automatically: limitation on LCG interface sites – global “switching-off” of LCG
 - Status descriptors for every stage of the job execution are essential for problem spotting and debugging
- Centralized and compact AliEn master services allow for fast updates and implementation of new features



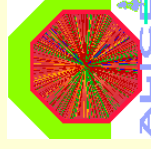
Considerations II

- LCG-2 provides many cycles
 - But required continuous efforts and interventions (ALICE and LCG)
 - Instabilities came in particular from the LCG site configurations
 - The LCG-SE is still very “fluid”, so we may expect instabilities
 - Analysis will be hopefully run with ARDA
- There is a big difference between pledged and available resources, and we learned this the hard way!

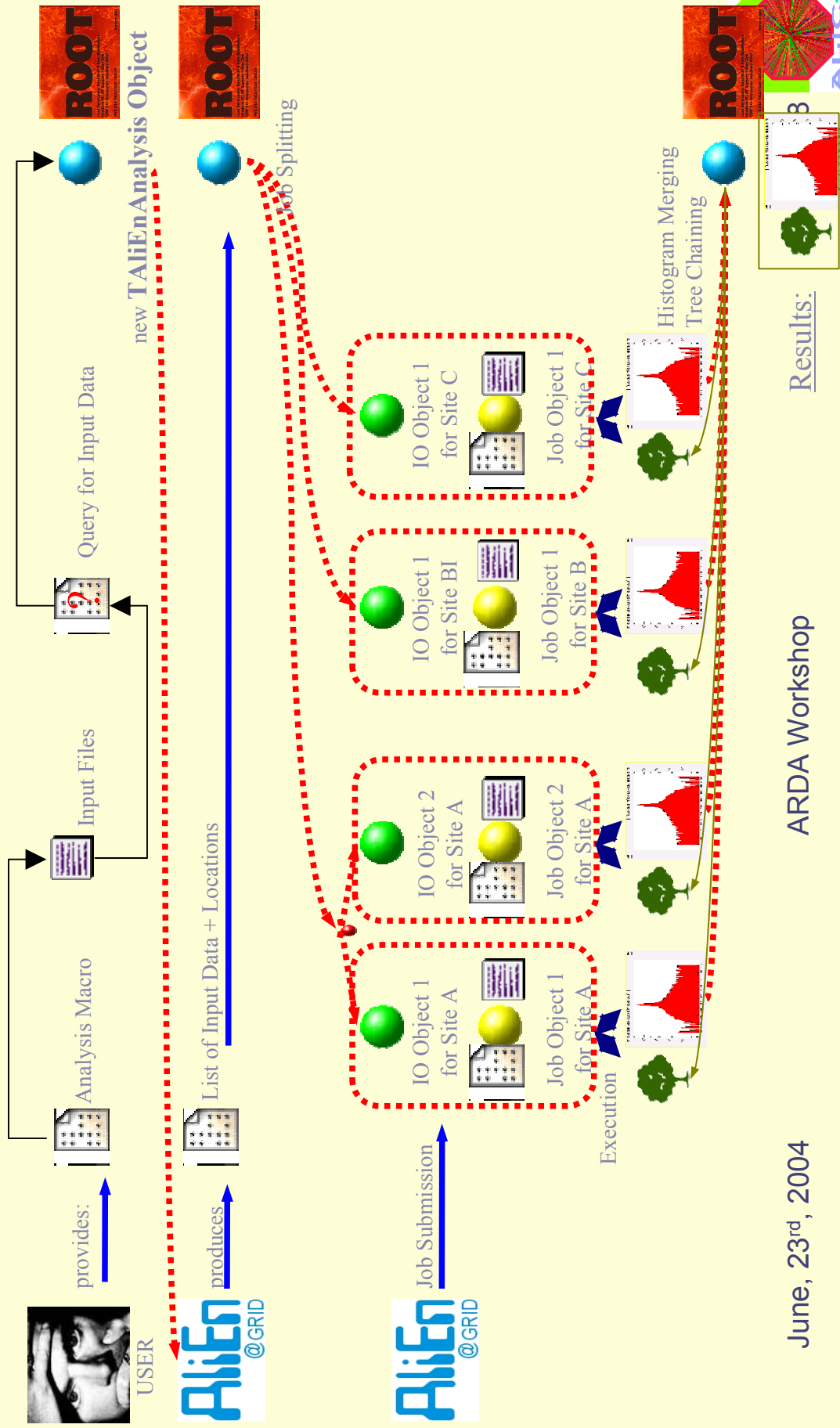


Considerations III

- The concept of AliEn as meta-grid works well, across three grids, and this is a success in itself
 - “keyhole” approach, some things become awkward (e.g. monitoring!)
- We are confident that we will reach the objectives

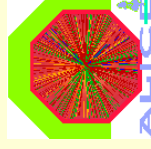


AliEn + ROOT / PROOF -> ARDA



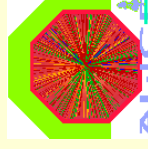
Considerations IV

- **Workload Management**
 - Limit to the amount of resources that can be served by a single “Master Queue”/“RB” – related to the job duration
 - Several Instances or a two-level hierarchy?
- **Data Management**
 - Minimize Data transfers (Application Data Model)
 - Functionality must be available, but the system must prevent misuse!
 - Association of a SE content with a DataCatalogue instance local to the SE, in addition to the Global “VO-Catalogue”?



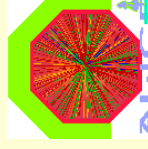
Considerations V

- A long term high-level requirement for the analysis: “local” Grid Services on a laptop?
 - Reproduce part of the Grid environment locally: development of “Grid-compliant” algorithms on a disconnected node
 - WM & DM short-circuit
 - Development of analysis algorithms, quick iterations with the same code/tools that will run on the Grid



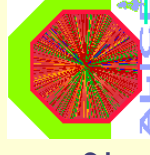
Considerations VI

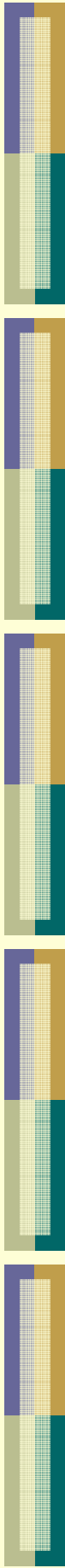
- ARDA
 - Multi-user Distributed Analysis **with** or **without** interactivity
 - AliEn/ROOT prototype available for batch mode
 - gLite/PROOF shall make it available in interactive mode
 - LCG-2 “as it is” not suitable for analysis
 - LCG-x backward compatibility is welcome until possible, but the evolution to gLite should give priority to the implementation of new functionality for interactive analysis in a WEB service-based design



Conclusions

- the DC is progressing with the help of all parties involved
- AliEn is fundamental, both as “Grid” and as “Meta-Grid” to LCG
- LCG support & performance is good
- LCG-2 shows some instabilities, but we are really pushing it!
- Heavy testing of the LCG-SE may bring more surprises, as large parts of it are new
- LCG-2 really looks and feels as a first-generation Grid
 - Very good for getting experience
 - Needing a new fresh start to build a production-grade system
- **Hope to test the first ARDA prototype in late summer/autumn**

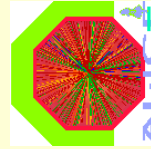




June, 23rd, 2004

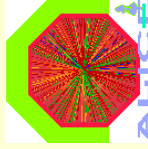
ARDA Workshop

33



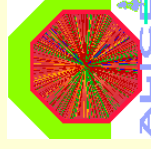
ALICE Physics Data Challenges

Period (<u>milestone</u>)	Fraction of the final capacity (%)	Physics Objective
<u>06/01-12/01</u>	1%	pp studies, reconstruction of TPC and ITS
<u>06/02-12/02</u>	5%	<ul style="list-style-type: none"> • First test of the complete chain from simulation to reconstruction for the PPR • Simple analysis tools • Digits in ROOT format
<u>01/04-06/04</u>	10%	<ul style="list-style-type: none"> • Complete chain used for trigger studies • Prototype of the analysis tools • Comparison with parameterised MonteCarlo • Simulated raw data
<u>05/05-07/05</u>	TBD	<ul style="list-style-type: none"> • Refinement of jet studies • Test of new infrastructure and MW • TBD
<u>01/06-06/06</u>	20%	<ul style="list-style-type: none"> • Test of the final system for reconstruction and analysis

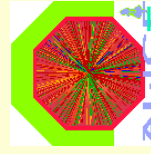
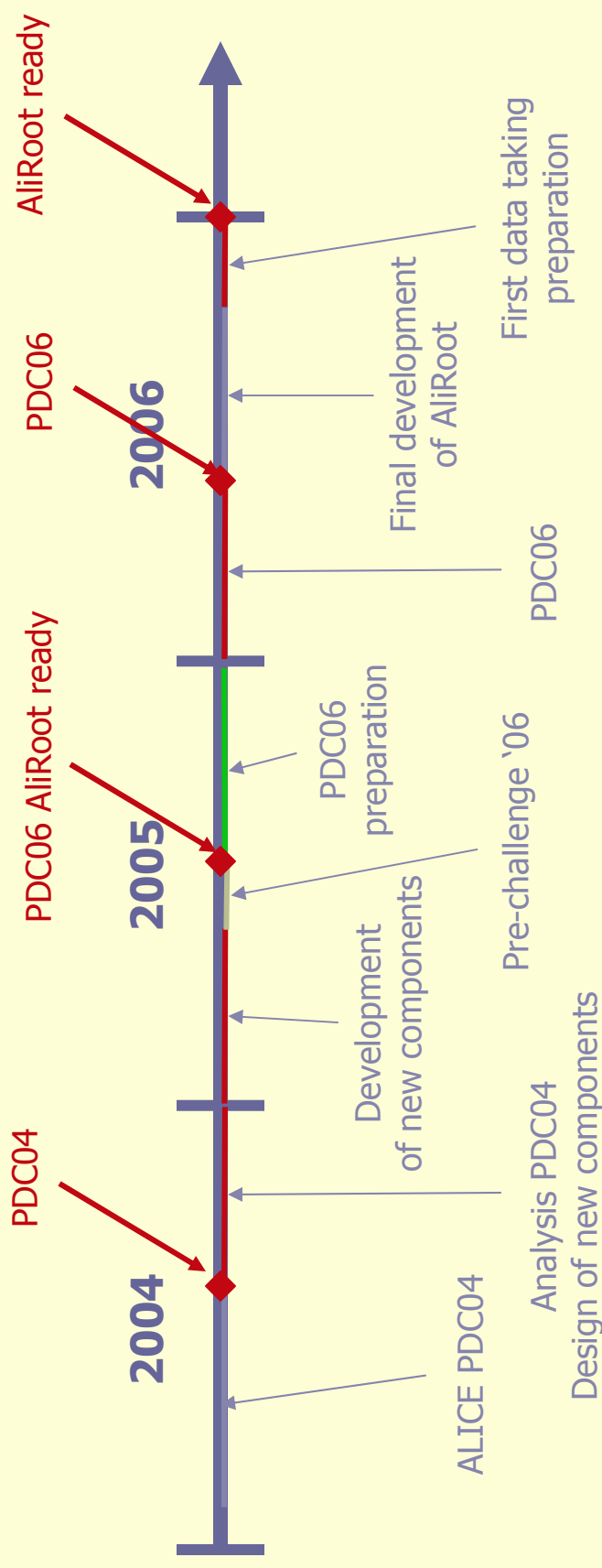


Why a new DC?

- We cannot stay 18 months without testing our “production capabilities”
- In particular we have to maintain the readiness of:
 - Code (AliRoot + MW)
 - ALICE distributed computing facilities
 - LCG infrastructure
 - Human “production machinery”
- Getting all the partners into “production mode” was a non-negligible effort
- We have to plan carefully size and physics objectives of this data challenge



ALICE Offline Timeline



Storage problems

- During most of last year 30TB of staging at CERN were in the plans for ALICE
 - This is what we requested
- When we started the DC we only had 2-3TB
- LCG has been very efficient in finding additional disks
- However we had to “throttle” production because otherwise we would have stopped anyway
- And we had to be very “creative” in swapping disk servers around

