



Enabling Grids for  
E-science in Europe

[www.eu-egee.org](http://www.eu-egee.org)

# NA4 Applications

F.Harris(Oxford/CERN)  
NA4/HEP coordinator



EGEE is a project funded by the European Union under contract IST-2003-508833

# Talk Outline

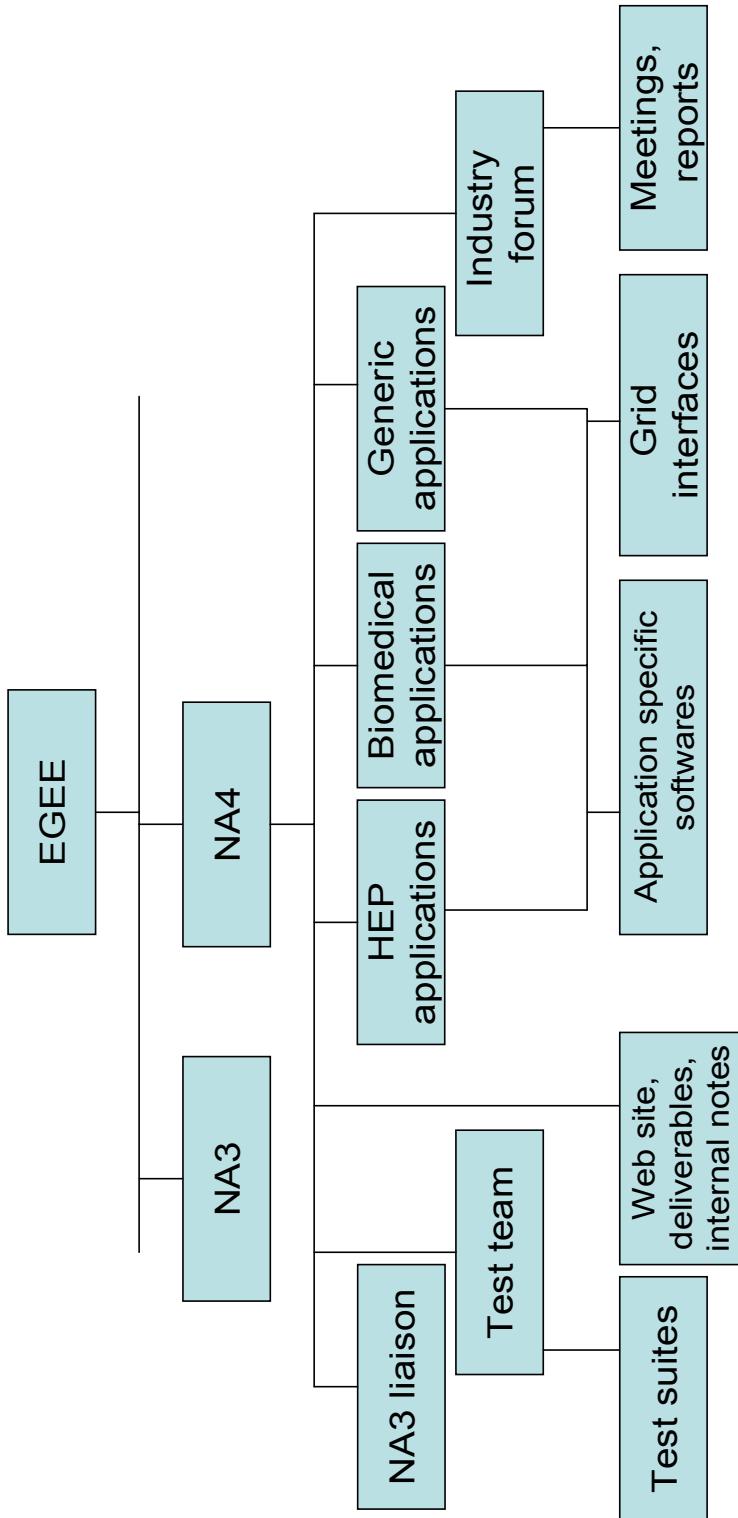
- The basic goals of NA4
- The organisation
- The participants and their roles
- The flavour of the work for the NA4 sub-groups
  - *Biomed*
  - *HEP*
  - ‘*Generic*’ *applications*
  - *Testing*
  - *Industry Forum*
- Milestones and deliverables
- Relations with other EGEE activities
- Concluding comments



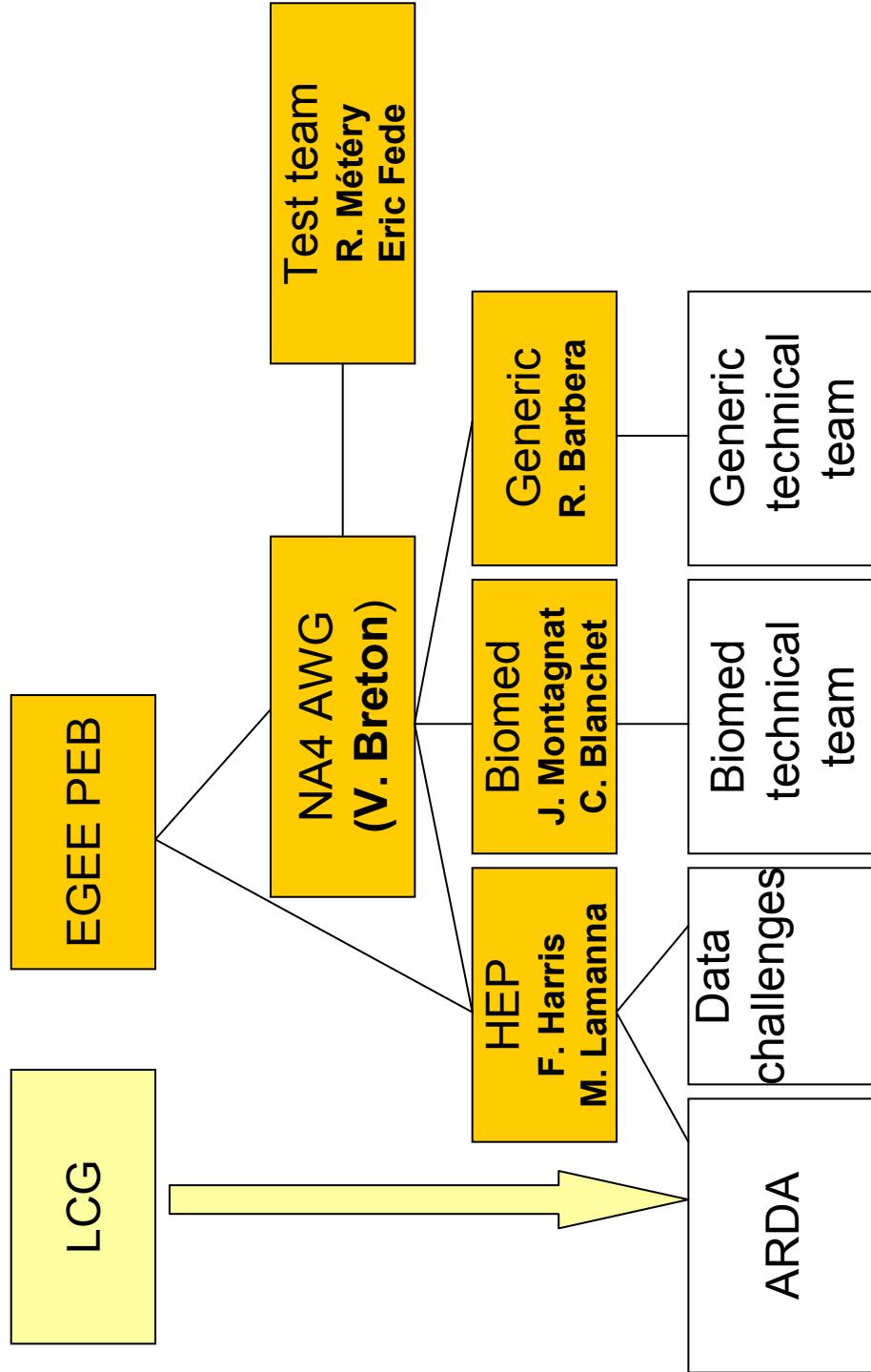
## NA4: Identification and support of early-user and established applications on the EGEE infrastructure

- To identify through the dissemination partners and a well defined integration process a portfolio of early user applications from a broad range of application sectors from academia, industry and commerce.
- To support development and production use of all of these applications on the EGEE infrastructure and thereby establish a strong user base on which to build a broad EGEE user community.
- To initially focus on two well-defined application areas – Particle Physics and Life sciences, while developing a process for supporting other application areas

# NA4 organisational structure



# NA4 technical organization



# Roles and staffing

<u>Federation</u>	<u>Role</u>	<u>FTE Funded</u>	<u>FTE Unfunded</u>
CERN	HEP Applications (coord.)	4	4 (9)
UK+Ireland	NA3 Liaison	0,5	0,5
Italy	Generic app (coord)	2	2
France	General coord., BioMed, Test team, Industry diss.	7	7
Northern Europe	Generic applications	1	1
Germany + Switzerland	Generic applications	1	1
Central Europe	Generic applications	1	1
South West Europe	BioMed	2	2
Russia	HEP, BioMed	3	3
<b>Totals</b>		<b>21,5</b>	<b>21,5</b>

# Biomedical Requirements

## ■ Complex data requirements

- ♦ Heterogeneous data formats (genomics, proteomics, image formats)
- ♦ Frequent data updates
- ♦ Complex data sets (medical records)
- ♦ Security/privacy constraints
- ♦ Long term archiving requirements

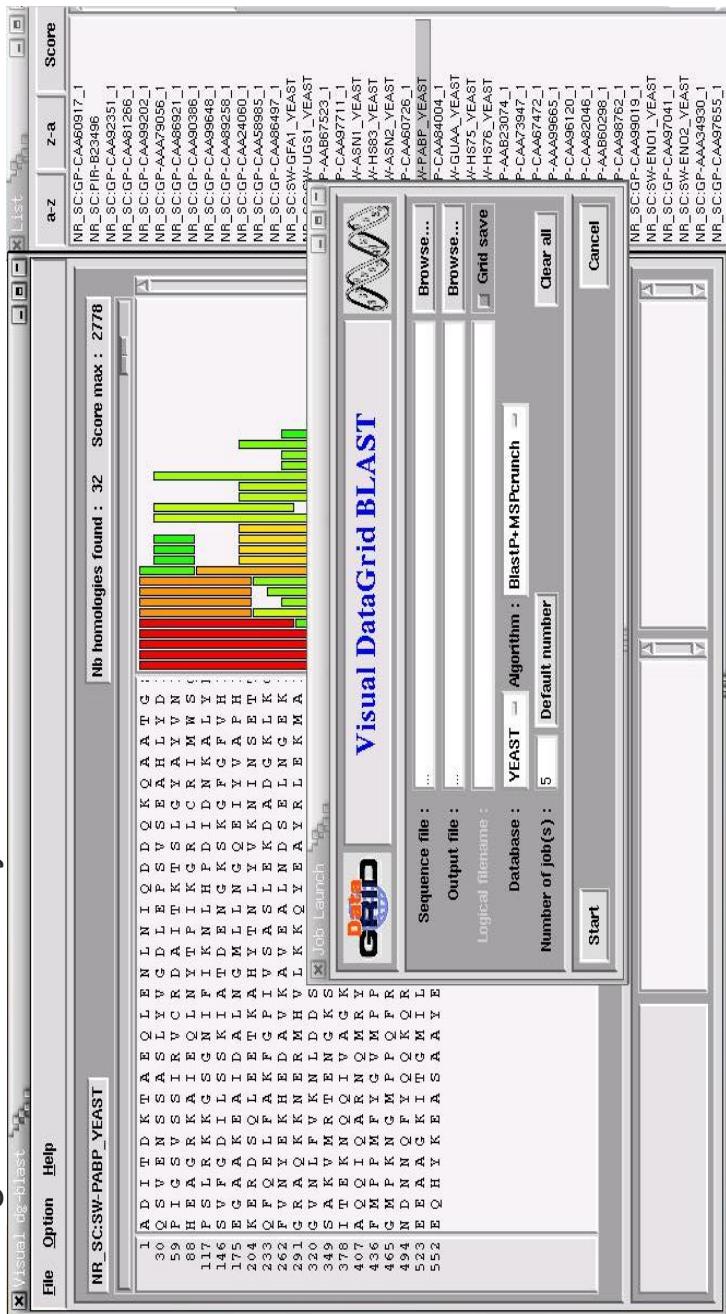
## ■ Complex processing requirements

- ♦ Bioinformatics: gene/proteome databases distributions
- ♦ Medical applications (screening, epidemiology...): image databases distribution
- ♦ Parallel algorithms for medical image processing, simulation, etc
- ♦ Interactive application (human supervision or simulation)
- ♦ Security/privacy constraints

# BLAST: Bioinformatics on the EDG testbed

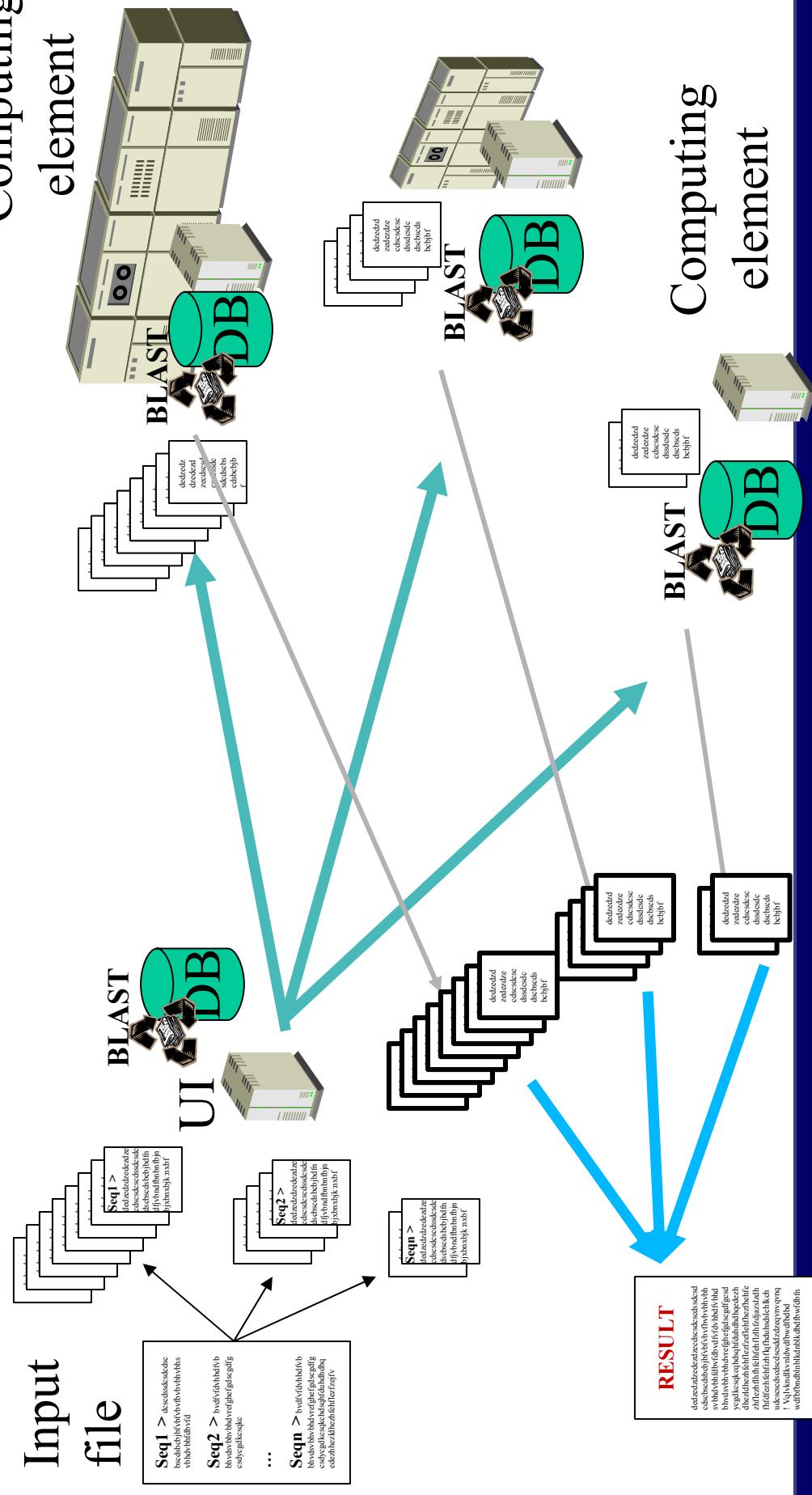
- BLAST is the first step for analysing new sequences: to compare DNA or protein sequences to other ones stored in personal or public databases. Ideal as a grid application.

- Requires resources to store databases and run algorithms
- Can compare one or several sequence against a database in parallel
- Large user community



# BLAST gridification

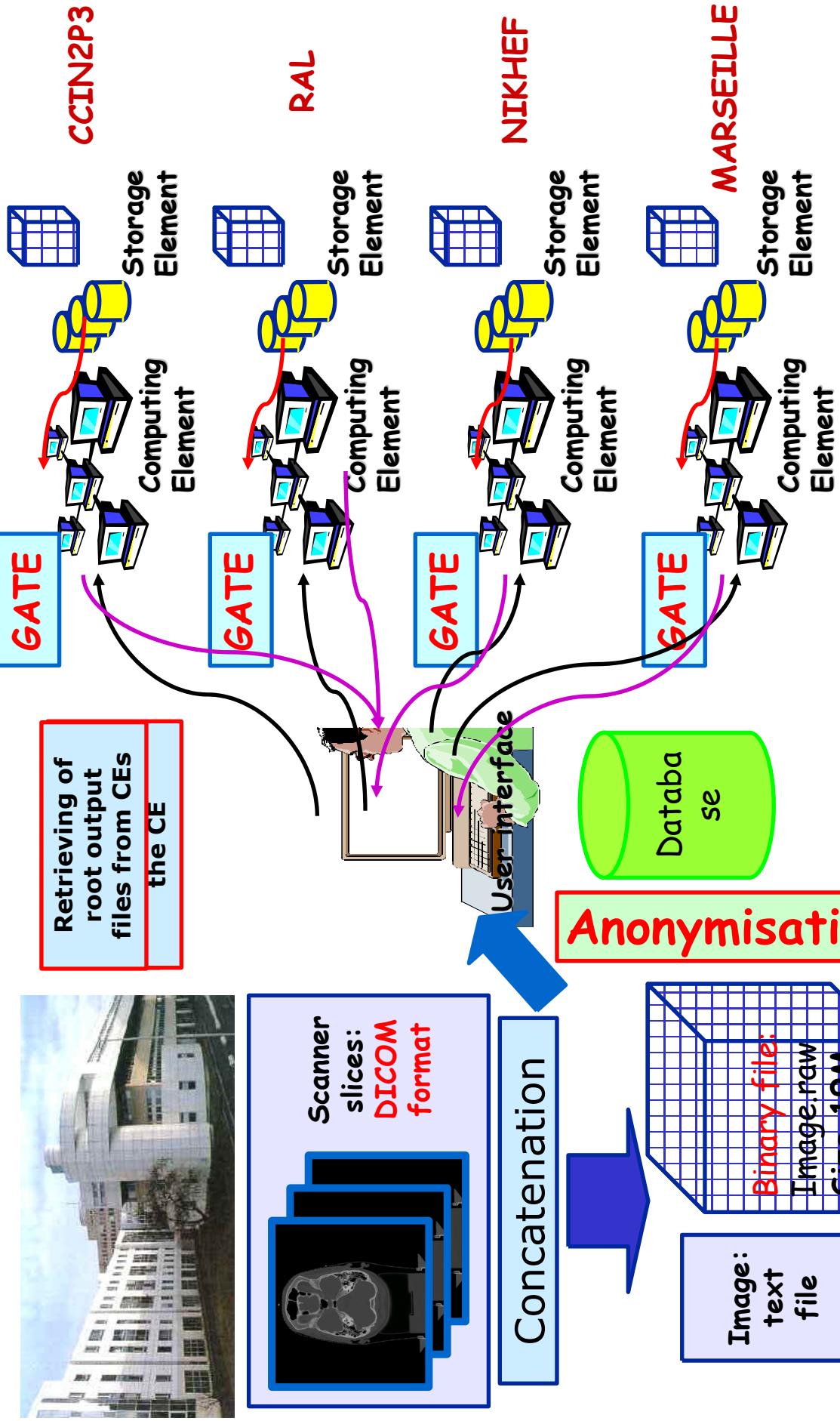
Computing element  
Input file



# Monte carlo simulation for radiotherapy planning



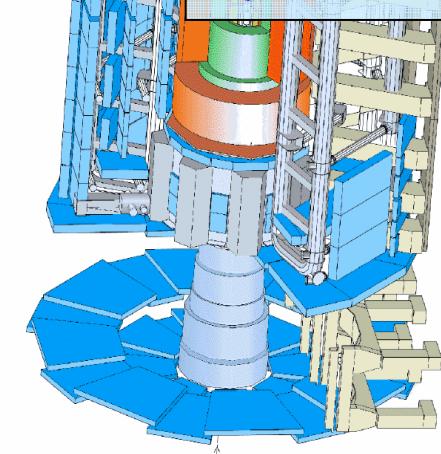
Enabling Grids for E-science in Europe



# LHC Experiments



ATLAS

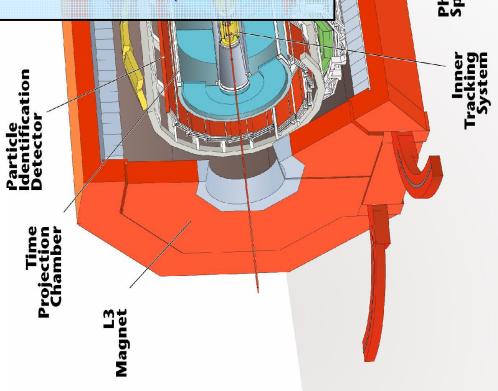


*Storage –*

Raw recording rate  $0.1 - 1 \text{ GByte/s}$

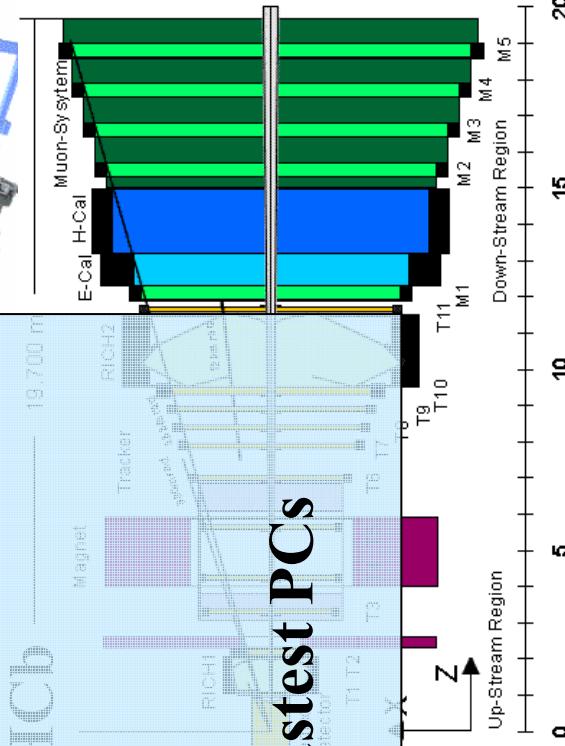
Accumulating at 5-8 PetaByte/year

ALICE



*Processing –*  
200,000 of today's fastest PCs

10 Petabyte of disk

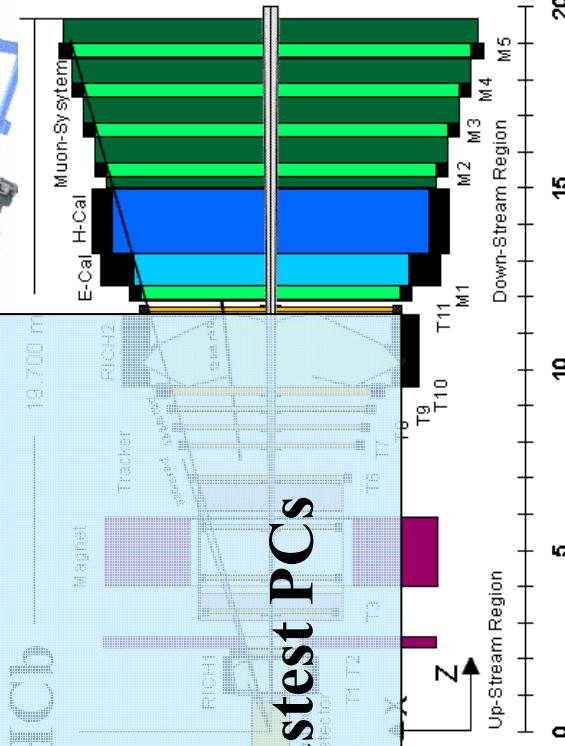


Raw recording rate  $0.1 - 1 \text{ GByte/s}$

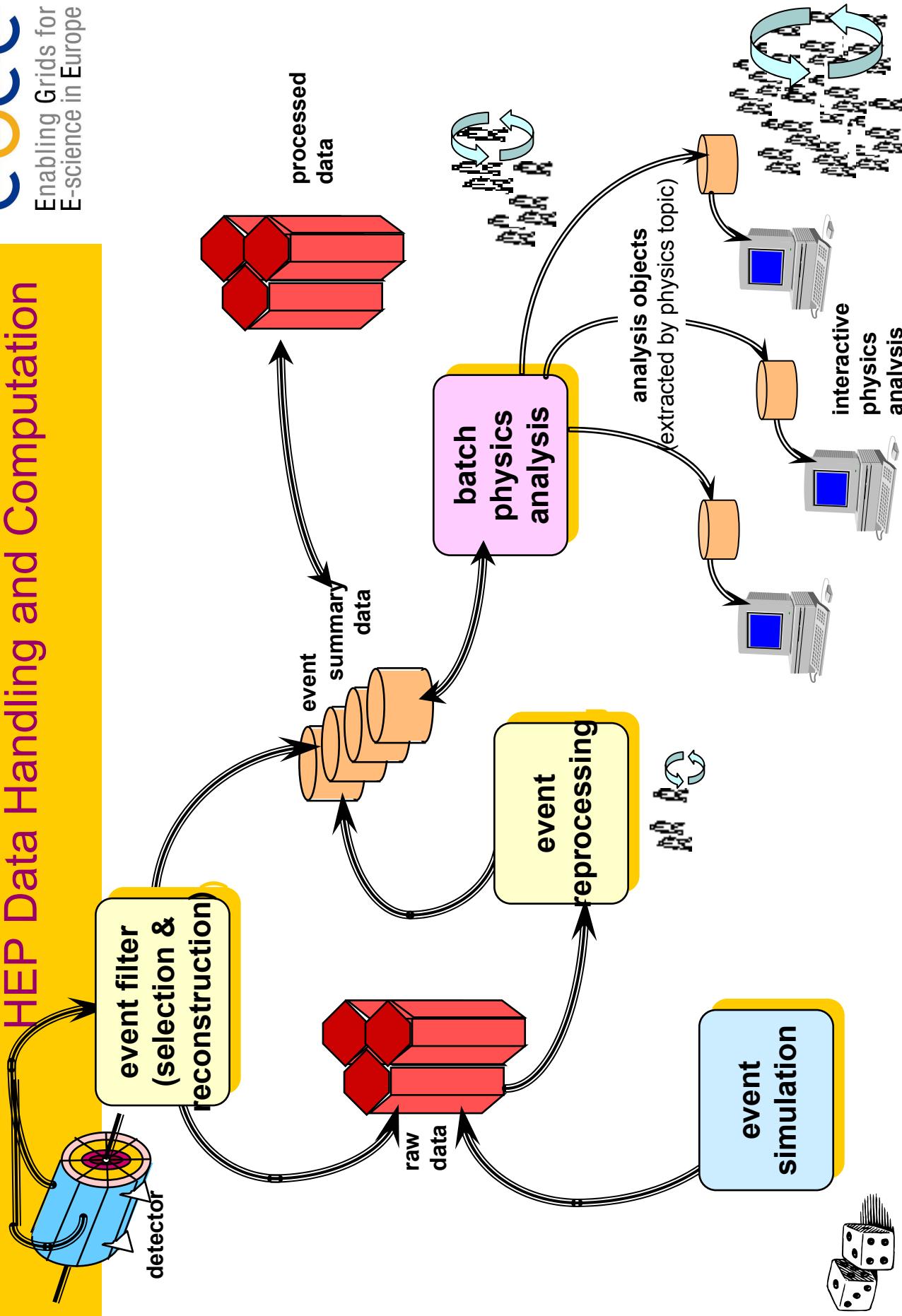
Accumulating at 5-8 PetaByte/year

*Processing –*  
200,000 of today's fastest PCs

10 Petabyte of disk

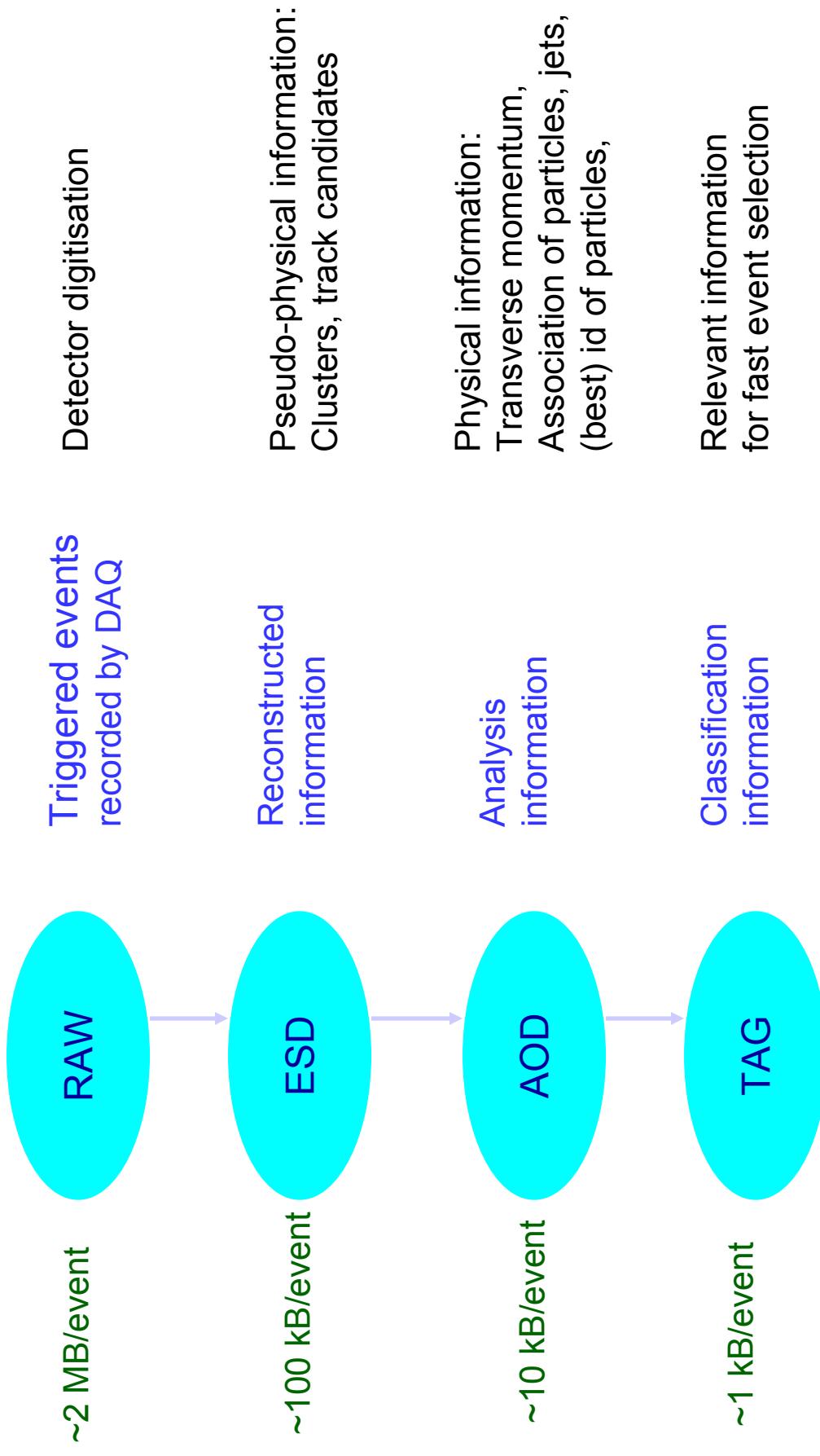


## HEP Data Handling and Computation



# Data Hierarchy

**“RAW, ESD, AOD, TAG”**



# HEP requirements

- **Data requirements**
  - ◆ Huge quantities of data( many Petabytes)
  - ◆ Data is of characteristic WORM (write once – read many times)
  - ◆ Data structured to allow data mining
  - ◆ Long time archiving, and need data copying around the world
- **Processing requirements**
  - ◆ Processing breaks down into 2 classes – production and ‘chaotic
    - ◆ Production done regularly with sets of ~  $10^{**9}$  events
    - ◆ Individual analyses done randomly on sets of maybe  $10^{**7}$  events – many hundreds of such users
  - ◆ Processing is highly parallel at ‘event’ level with directed graph like sequences in the processing
    - ◆ Interactive work very important for analysis- ability to save sessions and know origins of data (provenance)
    - ◆ Need global access to experiment databases for constants, running conditions etc.

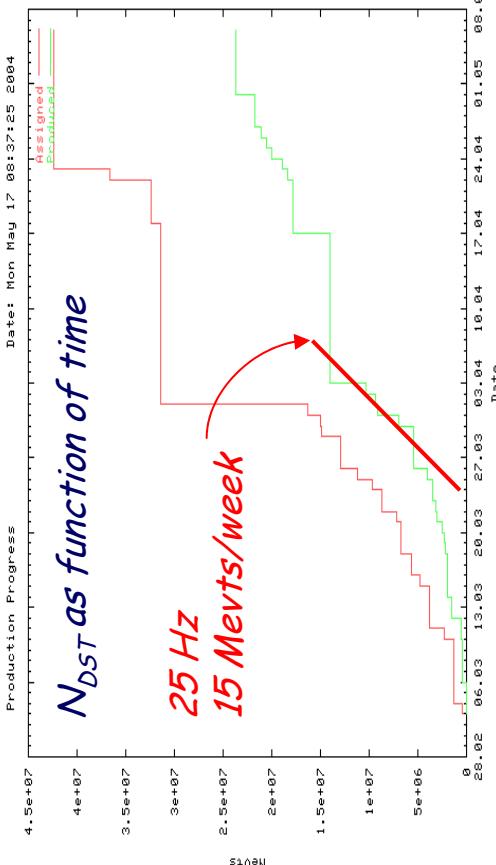
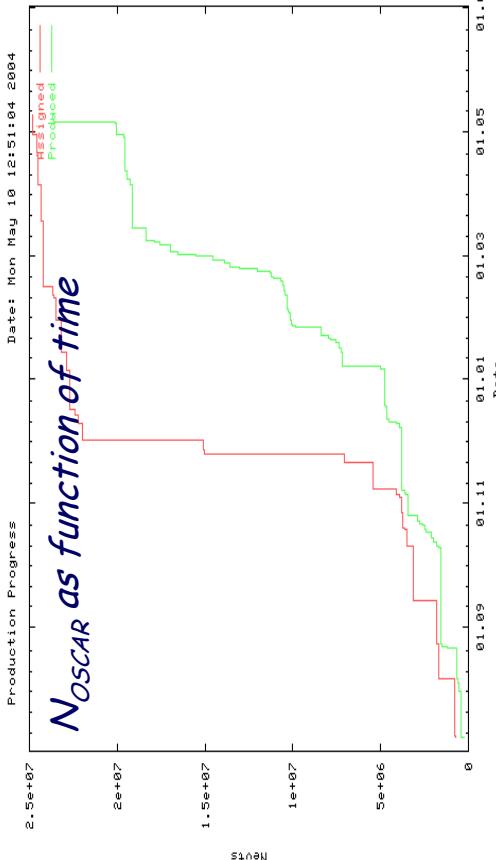
# Characteristics of CMS Data Challenge DC04

## Pre-Challenge Production

- Uses OCTOPUS
- After 8 months of continuous running:
  - 750,000 jobs
  - 3,500 KSI2000 months
  - 700,000 files
  - 80 TB of data
- With OSCAR (Geant 4)
  - 16 Mevts in 6 months

## Data Challenge

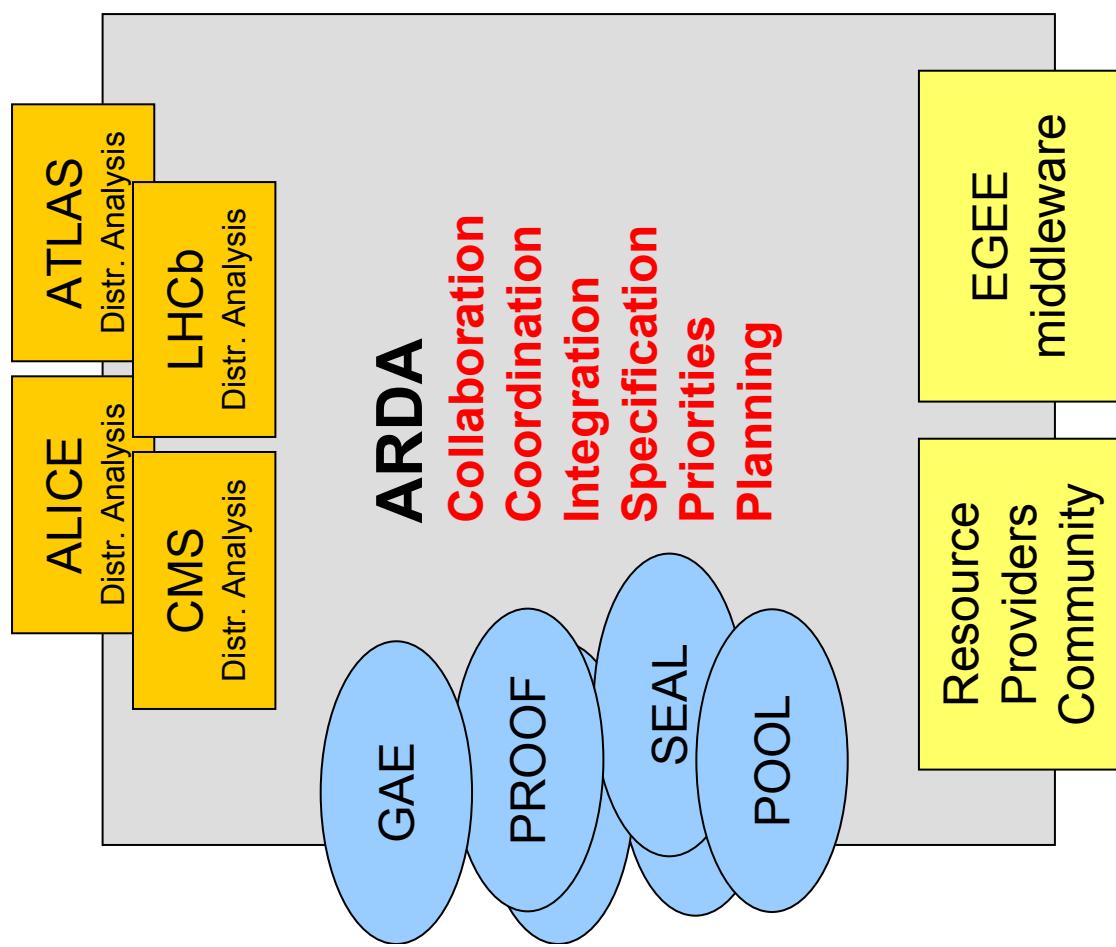
- Run the full data reconstruction and distribution chain at 25 Hz
- Achieved only for a limited amount of time:
  - 2,200 jobs/day (about 500 CPU's) running at Tier-0
  - 4 MB/s produced and distributed to each Tier-1
  - 0.4 files/s registered to RLS (with POOL metadata)



# ARDA Coordination and forum activities



Enabling Grids for  
E-science in Europe



## NA4 Generic Applications

- This is a key activity in the process of getting new scientific and industrial communities interested and committed to use the continental grid infrastructure built by the EGEE Project.
- GENIUS is a well established tool which will be fundamental in the process of interfacing new applications with the EGEE middleware hiding its complex internals to non-experts users from new communities.
- GILDA is a complete suite of grid elements (test-bed, CA, VO, monitoring system, web portal) and applications fully dedicated to dissemination purposes.  
**This could also represent the ideal grid testbed where to start the porting of new generic applications.**
- GILDA is the dissemination tool which will be used by NA3 during courses and tutorials so the important aspect of induction of the grid paradigm to new communities is also covered.
- It is now important to have the first meeting of the EGAAAP board and define the first Generic Applications to be interfaced.

# The GILDA home page (<http://gilda.ct.infn.it>)



## GILDA ( Grid INFN Laboratory for Dissemination Activities )

Grid tutorials  
Instructions for users  
Instructions for sites  
Useful links  
Usage Statistics

GILDA consists of the following elements:

- the GILDA Testbed: a series of sites spread all over Italy where the last version of the GridIt grid middle-ware is installed;
- the GILDA Certification Authority: a fully functional Certification Authority which issues 14-days X.509 certificates to everybody wanting to experience grid computing on the GILDA Testbed;
- the GILDA Virtual Organization: a Virtual Organization gathering all people wanting to experience grid computing on the GILDA Testbed;
- the Grid Demonstrator: a customized version of the full GENIUS web portal, jointly developed by INFN and NICE, from where users belonging to the GILDA VO can submit a pre-defined set of applications to the GILDA Testbed;
- the GENIUS web portal: the full GENIUS web portal, to be used only during grid tutorials;
- the monitoring system: a versatile monitoring system completely based on GridICE, the grid monitoring tool developed by INFN;
- the GILDA mailing list: gilda@infn.it, also archived on the web [here](#).

GILDA is an activity of the Italian Istituto Nazionale di Fisica Nucleare (INFN) carried on in the context of both the Italian [INFN Grid](#) and European [EGEE](#) Projects.



# The Generic Application questionnaire (contribution to MNA4.1)



Enabling Grids for  
E-science in Europe

- Questionnaire to get information and first requirements from new communities interested in using the EGEE Infrastructure (<http://alipc1.ct.infn.it/grid/egee/na4/questionnaire/na4-genapp-questionnaire.doc>)
  - Feed-backs received so far
    - (<http://alipc1.ct.infn.it/grid/egee/na4/questionnaire>):
  - **Astrophysics (EVO and Planck satellite)**
  - **Earth Observation (ozone maps, seismology, climate)**
  - **Digital Libraries (DILIGENT Project)**
  - **Grid Search Engines (GRACE Project)**
  - **Industrial applications (SIMDAT Project)**
- Interest also from Computational Chemistry (Italy and Czech Republic), Civil Engineering (Spain), and Geophysics (Switzerland and France) communities

# EGEE Industry Forum Objectives

- The main role of the Industry Forum in the Enabling Grids for E-Science in Europe (EGEE) project is to raise awareness of the project amongst industry and to encourage businesses to participate in the project. This will be achieved by making direct contact with industry, liaising between the project partners and industry in order to ensure businesses get what they need in terms of standards and functionality and ensuring that the project benefits from the practical experience of businesses.
- The members of the EGEE Industry Forum are companies of all sizes who have their core business or a part of their business within Europe.
- The Industry Forum will be managed by a steering group consisting of EGEE project partners and representatives from business
  - <http://public.eu-egee.org/industry-forum/information>

# NA4 Testing Group



**Three types of tests will be developed:**

(based on user requirements and on the experience gathered by the ongoing activities like LHC DCs and ARDA prototyping)

- **Tests of service availability:** This set of tests will check the EGEE services availability. All the services providing by the GRID should be tested : Job submission and management, files management, Information service, ...
- **Tests of functionality:** To verify that the functionalities required are available , usable and complete; for example, file creation, moving and deletion, information publication, errors recovery.
- **Tests to measure performances:** Their goal is to characterize the testbed from the end users/application perspective. Part of them will be time measurements ( time to submit X job, time to replicate Y files,...), others will address scalability measurements ( how many jobs can be accepted by service Z, files limits size,...) while others will be more abstract (information availability, errors message access,...).
- **This work should be done in close collaboration with ARDA , JRA1 and SA1**

# Milestones for NA4 applications



First applications migrated to the EGEE infrastructure		
MNA4.1	M6	<ul style="list-style-type: none"><li>•HEP data challenges for 4 LHC experiments and for D0</li><li>•Biomedicine – GATE simulation in nuclear medicine + others</li><li>•Plus the first ‘generic’ applications</li></ul>
MNA4.2	M12	<p>First <b>external review</b> of Applications Identification and Support with feedback</p>
MNA4.3	M24	<ul style="list-style-type: none"><li>•<b>Second external review</b> of Applications Identification and Support with feedback</li></ul>

# Deliverables for NA4 applications



DNA4.1	M3	<b>Definition of Common Application Interface</b> (especially important for new applications...probably based on GENIUS flavoured work)
DNA4.2	M6	<b>Target Application Sector Strategy document</b> (work for getting new applications onto EGEE)
DNA4.3	M9	<b>EGEE Application Migration Progress report</b> (revision M15 and M21) •All applications will include evaluations on production and pre-production services of LCG (current and 'new' middleware)
DNA4.4	M24	<b>Final Report of Application Identification and Support Activity</b>

# NA4 relations with other EGEE activities and other bodies (1)

- **SA1 grid operations**

- How to get new VOs onto LCG from different domains?
- How to integrate new resources(sites) into LCG coming from different application areas?
- Rationalisation of test procedures
- Working with national agencies (e.g. GridPP Application monitoring)

- **NA3 training**

- Estimating requirements for courses
- Design and implementation of courses

- **JRA1 middleware**

- All applications input requirements and monitor their satisfaction with feedback to middleware (process goes through the PTF-Project Technical Forum)

- **JRA2 quality assurance**

- NA4 have a representative on this group to define process for monitoring quality of EGEE services

## NA4 relations with other EGEE activities and other bodies (2)

- **JRA3 security**
  - Security of medical (and other application) data
  - Security for sites
- **SA2,JRA4 networking**
  - Global HEP requirements through LCG
  - Biomed and other applications must similarly give global needs
  - NA4 will give individual application use cases especially where problems have been encountered
- **LCG**
  - NA4/HEP are presented on the LCG/GAG(grid Applications Group)
    - This is HEP source of requirements and giving feedback to middleware on ‘customer satisfaction’. Some GAG people are on the PTF.

# Conclusions

- **NA4 is up and running now**
  - HEP is using LCG-2 for data challenges
  - ARDA is well under way and waiting for first new middleware prototype
  - Biomedicine has applications ready to go onto LCG-2 and pre-production services
  - Generic group is very active with GILDA and excellent relations with NA3
  - Testing group has active dialogue with JRA1 and ARDA for rationalising testing effort
  - Industry forum has developed links with several companies (see EGEE Cork presentations)
- **NA4 open meeting Jul 14-16 at Catania with emphasis on inter-activity dialogue** (with middleware,operations,security,networking)
  - **NA4 Web site** <http://egee-na4.ct.infn.it>