



Introduction to GRID computing and overview of the European Data Grid Project



The European DataGrid Project

<http://www.edg.org>



Overview

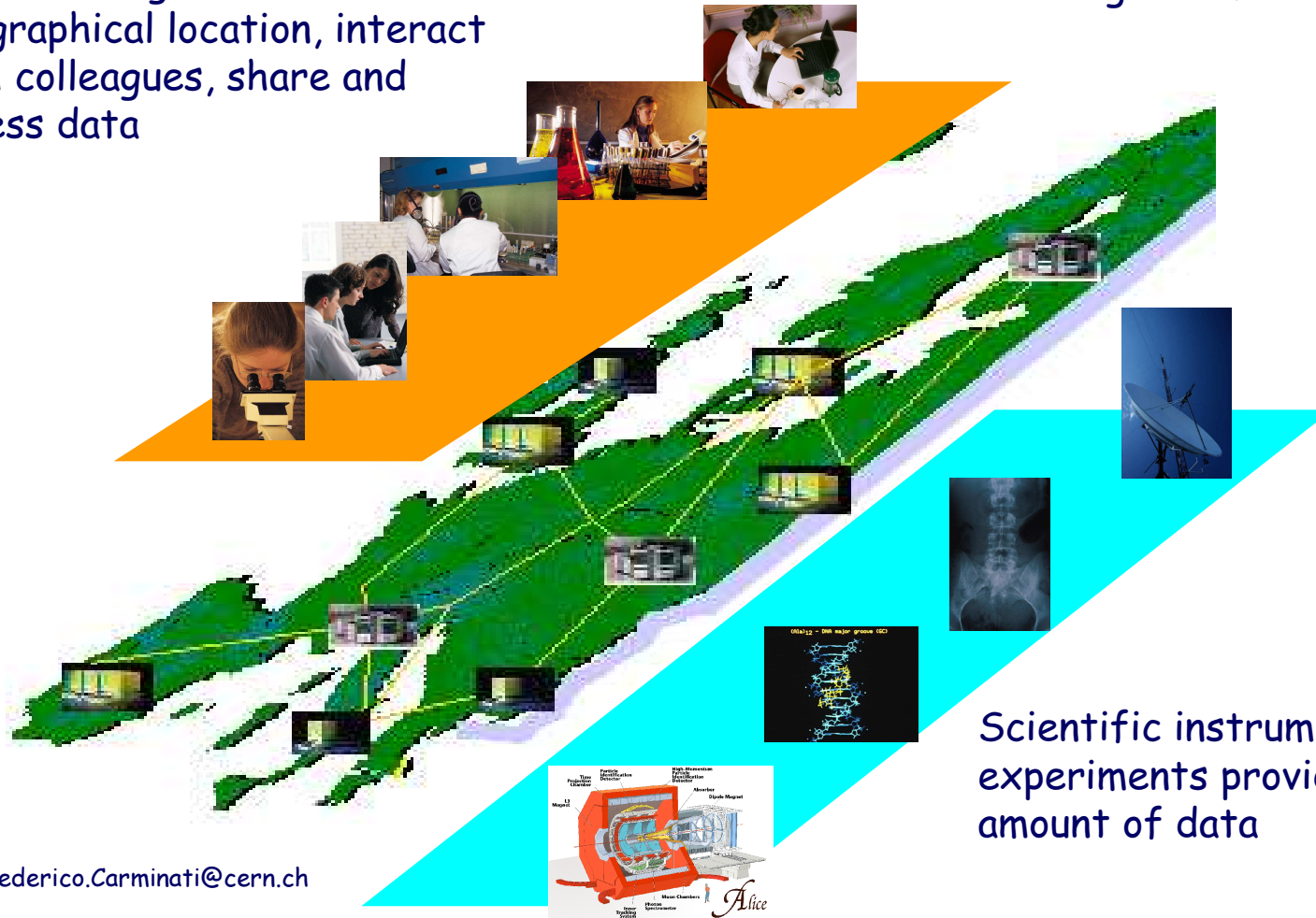
- What is GRID computing ?
- What is a GRID ?
- Why GRIDs ?
- GRID projects world wide
- The European Data Grid
 - Overview of EDG goals and organization
 - Overview of the EDG middleware components



The Grid Vision

Researchers perform their activities regardless geographical location, interact with colleagues, share and access data

The GRID: networked data processing centres and "middleware" software as the "glue" of resources.



Scientific instruments and experiments provide huge amount of data

Federico.Carminati@cern.ch



What is GRID computing :

➤ **coordinated resource sharing and problem solving in dynamic, multi-institutional virtual organizations.** [I.Foster]

- A VO is a collection of users sharing similar needs and requirements in their access to processing, data and distributed resources and pursuing similar goals.

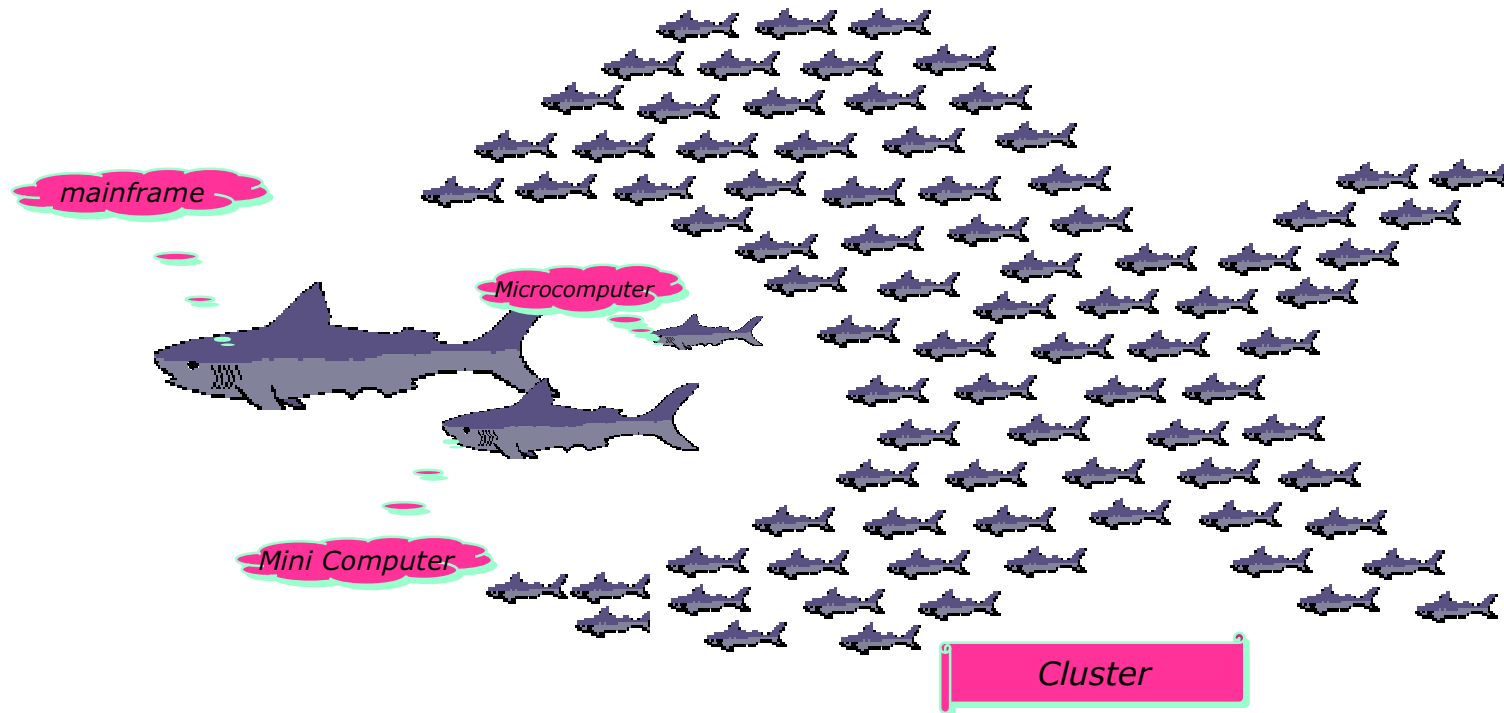
➤ **Key concept :**

- **ability to negotiate resource-sharing arrangements among a set of participating parties (providers and consumers) and then to use the resulting resource pool for some purpose** [I.Foster]



The GRID distributed computing idea 1/2

Once upon a time.....

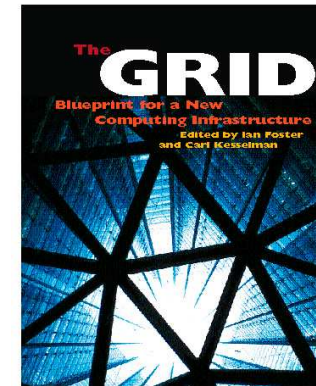
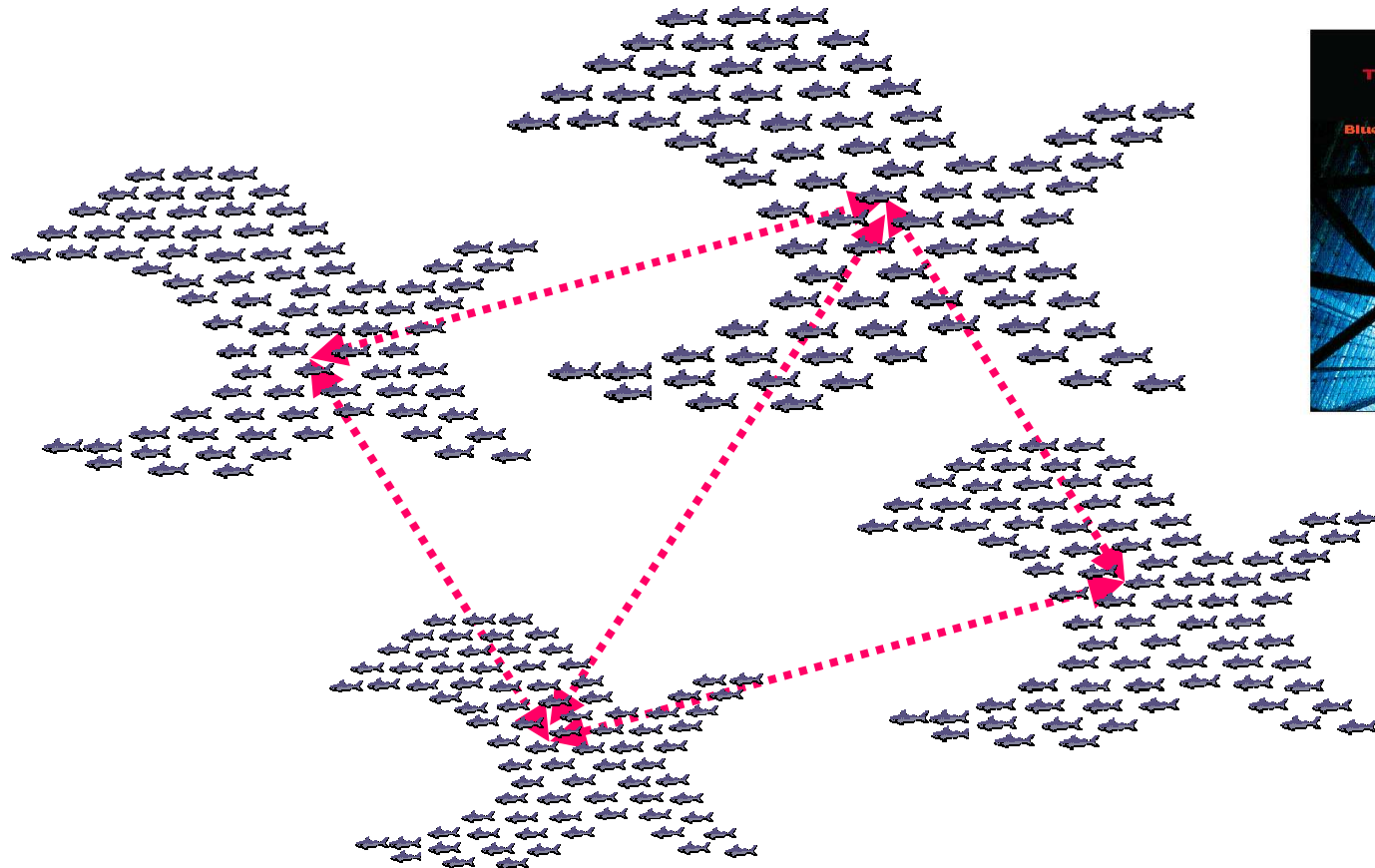


(by Christophe Jacquet, after N.O.W.)



The GRID distributed computing idea 2/2

...and today



(by Christophe Jacquet)



Differences between grids and distributed applications

- Distributed applications already exist, but they tend to be specialised systems intended for a single purpose or user group
- Grids go further and take into account:
 - Different kinds of resources
 - Not always the same hardware, data and applications
 - Different kinds of interaction
 - User groups or applications want to interact with grids in different ways
 - Dynamic nature
 - Resources and User Groups added/removed/changed frequently



Main characteristics of a grid architecture

- Service providers
 - Publish the availability of their services via information systems
 - Such services may *come-and-go or change* dynamically
 - E.g. a testbed site that offers x CPUs and y GB of storage
- Service brokers
 - Register and categorize published services and provide search capabilities
 - E.g. 1) EDG Resource Broker selects the best site for a "job"
2) Catalogues of data held at each testbed site
- Service requesters
 - Single sign-on: log into the grid once
 - Use brokering services to find a needed service and employ it
 - E.g. CMS physicists submit a simulation job that needs 12 CPUs for 6 hours and 15 GB which gets scheduled, via the Resource Broker, on the CERN testbed site



GRID security

- Resource providers are essentially "opening themselves up" to itinerant users
- Secure access to resources is required
 - X.509 Public Key Infrastructure
- User's identity has to be certified by (mutually recognized) national Certification Authorities (CAs)
- Resources (node machines) have to be certified by CAs
- Temporary delegation from users to processes to be executed "in user's name" (proxy certificates)
- Common agreed policies for accessing resource and handling user's rights across different domains within Virtual Organizations



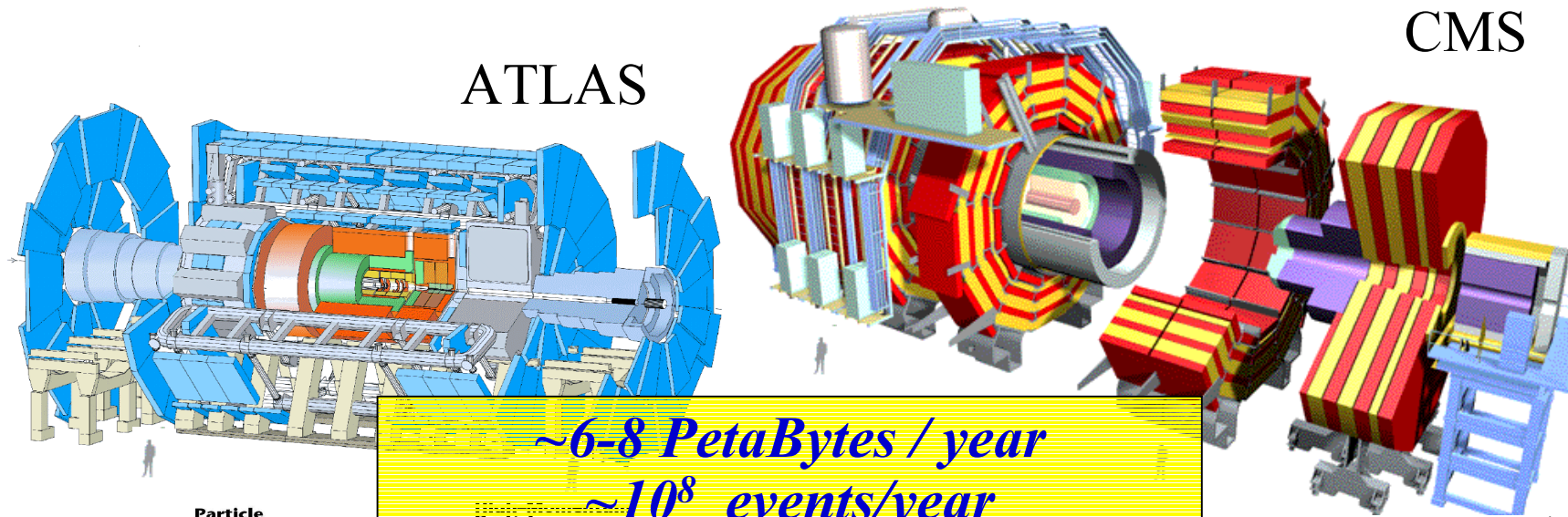
Why GRIDs

- **Scale of the problems**
 - frontier research in many different fields today requires world-wide collaborations (i.e. multi-domain access to distributed resources)
- **GRIDs provide access to large data processing power and huge data storage possibilities**
 - As the grid grows its usefulness increases (more resources available)
- **Large communities of possible GRID users :**
 - High Energy Physics
 - Environmental studies: Earthquakes forecast, geologic and climate changes, ozone monitoring
 - Biology, Genetics, Earth Observation
 - Astrophysics,
 - New composite materials research
 - Astronautics, etc.

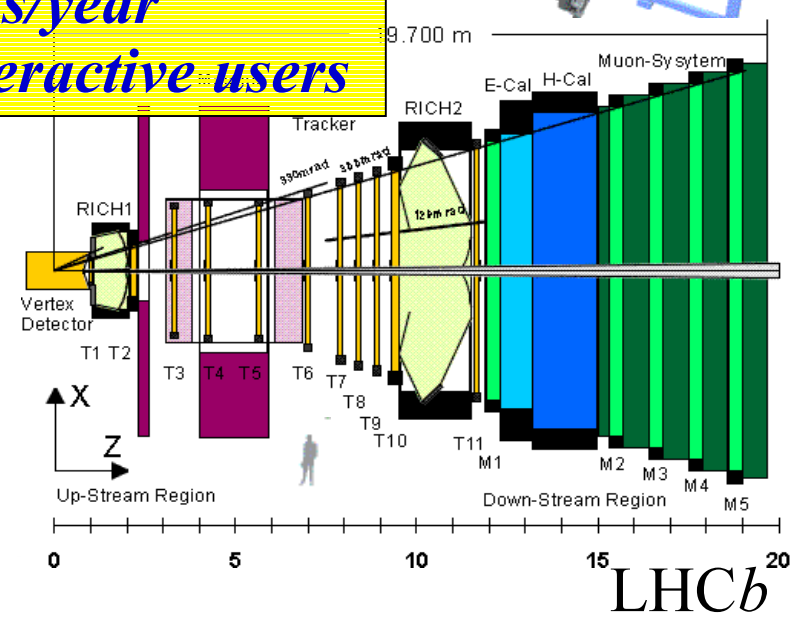
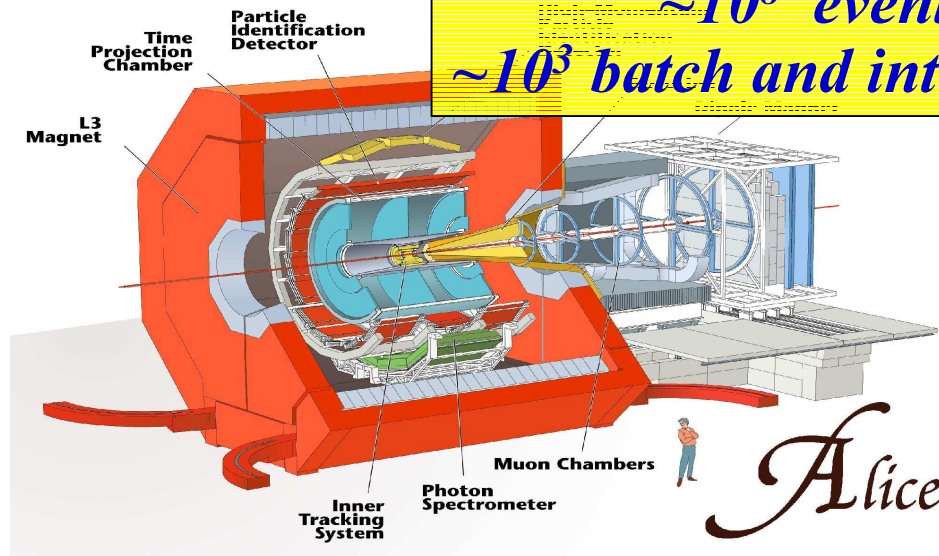


High Energy Physics

The LHC Detectors



~6-8 PetaBytes / year
~10⁸ events/year
~10³ batch and interactive users



Federico.carminati , EU review presentation

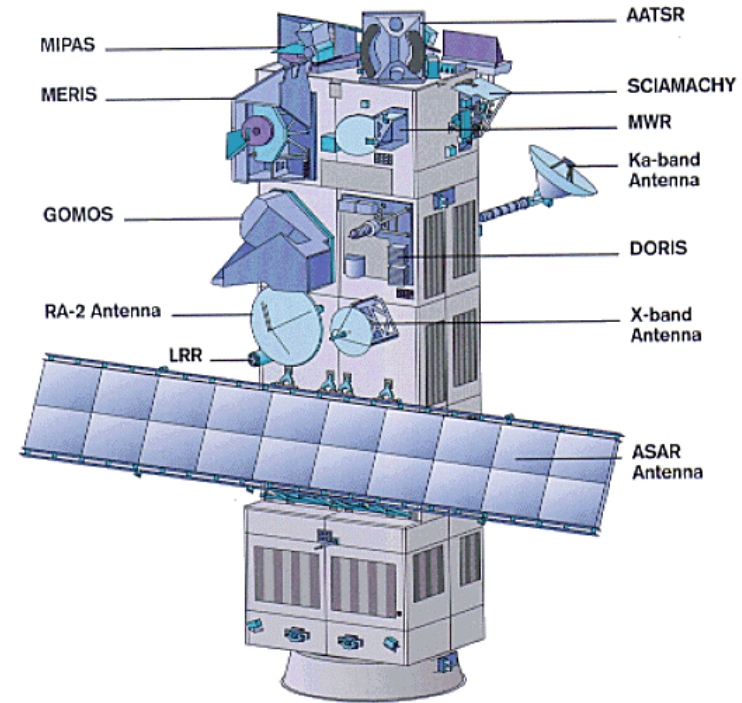
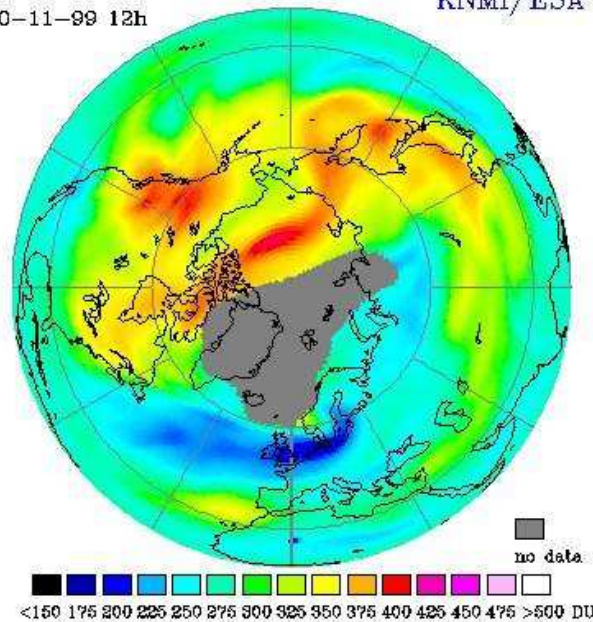


Earth Observation

ESA missions:

- about 100 Gbytes of data per day (ERS 1/2)
- 500 Gbytes, for the next ENVISAT mission (2002).

Assimilated GOME total ozone
30-11-99 12h
KNMI/ESA



DataGrid contribute to EO:

- enhance the ability to access high level products
- allow reprocessing of large historical archives
- improve Earth science complex applications (data fusion, data mining, modelling ...)

Source: L. Fusco, June 2001

The screenshot displays the Visual de-blast application window. The main window shows a sequence alignment for 'NR_SC:SW-PABP_YEAST' with 32 homologues found and a maximum score of 2778. A bar chart visualizes the alignment scores. A 'Job Launch' dialog box is open in the foreground, titled 'Visual DataGrid BLAST'. The dialog includes fields for 'Sequence file', 'Output file', and 'Logical filename', each with a 'Browse...' button. It also features a 'Grid save' checkbox, a 'Database' dropdown set to 'YEAST', an 'Algorithm' dropdown set to 'BlastP+MSFcrunch', and a 'Number of job(s)' field set to '5'. The dialog has 'Start' and 'Cancel' buttons.

Visual de-blast
 File Option Help

NR_SC:SW-PABP_YEAST Nb homologues found : 32 Score max : 2778

```

1  A D I T D K T A E Q L E N L N I Q D D Q K Q A A T G :
30 Q S V E N S S A S L Y V C D L E P S V S E A H L Y D :
59 P I G S V S S I R V C R D A I T K T S L G Y A Y V N :
88 H E A C R K A I E Q L N Y T P I K C R L C R I M W S :
117 P S L R K K G S G N I F I K N L H P D I D N K A L Y :
146 S V F G D I L S S K I A T D E N G K S K C F C F V H :
175 E G A A K E A I D A L N G M L L N G Q E I Y V A P H :
204 K E R D S Q L E E T K A H Y T N L Y V K N I N S E T :
233 Q F Q E L F A K F G P I V S A S L E K D A D G K L K :
262 F V N Y E K H E D A V K A V E A L N D S E L N G E K :
291 C R A Q K K N E R M H V L K K Q Y E A Y R L E K M A :
320 G V N L F V K N L D D S :
349 S A K V M R T E N G K S :
378 I T E K N Q Q I V A G K :
407 A Q Q I Q A R N Q M R Y :
436 F H P P M F Y G V M P P :
465 C H P K N G M F P Q P R :
494 N D N N Q F Y Q Q K Q R :
523 E E A A G K I T G M I L :
552 E Q H Y K E A S A A Y E
  
```

Visual DataGrid BLAST

Sequence file : Browse...

Output file : Browse...

Logical filename : Grid save

Database : YEAST Algorithm : BlastP+MSFcrunch

Number of job(s) : 5 Default number Clear all

Start Cancel

List

a-z	Z-a	Score
NR_SC:GP-CAA60917_1		
NR_SC:PIR-B23496		
NR_SC:GP-CAA82351_1		
NR_SC:GP-CAA81266_1		
NR_SC:GP-CAA99202_1		
NR_SC:GP-AAA79056_1		
NR_SC:GP-CAA86921_1		
NR_SC:GP-CAA80386_1		
NR_SC:GP-CAA99648_1		
NR_SC:GP-CAA89258_1		
NR_SC:GP-CAA24060_1		
NR_SC:GP-CAA58985_1		
NR_SC:GP-CAA86497_1		
NR_SC:SW-GFA1_YEAST		
NR_SC:SW-UGS1_YEAST		
P-AA867523_1		
P-CAA97711_1		
W-ASN1_YEAST		
W-HS83_YEAST		
W-ASN2_YEAST		
P-CAA80726_1		
W-PABP_YEAST		
P-CAA84004_1		
W-GLAA_YEAST		
W-HS75_YEAST		
W-HS76_YEAST		
P-AA823074_1		
P-CAA73947_1		
P-CAA67472_1		
P-AAA99685_1		
P-CAA96120_1		
P-CAA82046_1		
P-AA860298_1		
P-CAA86762_1		
NR_SC:GP-CAA99019_1		
NR_SC:SW-ENO1_YEAST		
NR_SC:GP-CAA97041_1		
NR_SC:SW-ENO2_YEAST		
NR_SC:GP-AAA34930_1		
NR_SC:GP-CAA97655_1		



GRID projects world wide

➤ EU

- EDG (EU-IST) - R&D EU GRID project [www.edg.org]
- CrossGRID - QoS - interactive apps. [www.crossgrid.org]
- DataTAG - inter-operability (EU-USA) [www.datatag.org]
- LCG - The LHC Computing GRID - Deployment [cern.ch/lcg]
- The new 16,2 B Euro EU VI Framework Prog. GEANT based GRID projects

➤ USA

- GriPhyN [www.griphyn.org]
- iVDGL-VDTv1 [www.idvgl.org]
- PPDG (NSF, DoE) [www.ppdg.org]

➤ Asia

- ApGrid [www.apgrid.org]
- Pragma (USA-Asia)

And many more. . .



The European Data Grid

- To build on the emerging Grid technology to develop a sustainable computing model for effective share of computing resources and data

➤ Start : Jan 1, 2001 End : Dec 31, 2003

- Specific project objectives:
 - Middleware for fabric & Grid management (mostly funded by the EU)
 - Large scale testbed (mostly funded by the partners)
 - Production quality demonstrations (partially funded by the EU)
- To collaborate with and complement other European and US projects
- Contribute to Open Standards and international bodies:
 - Co-founder of Global GRID Forum and host of GGF1 and GGF3
 - Industry and Research Forum for dissemination of project results



The EDG Main Partners

- CERN - International (Switzerland/France)
- CNRS - France
- ESA/ESRIN - International (Italy)
- INFN - Italy
- NIKHEF - The Netherlands
- PPARC - UK





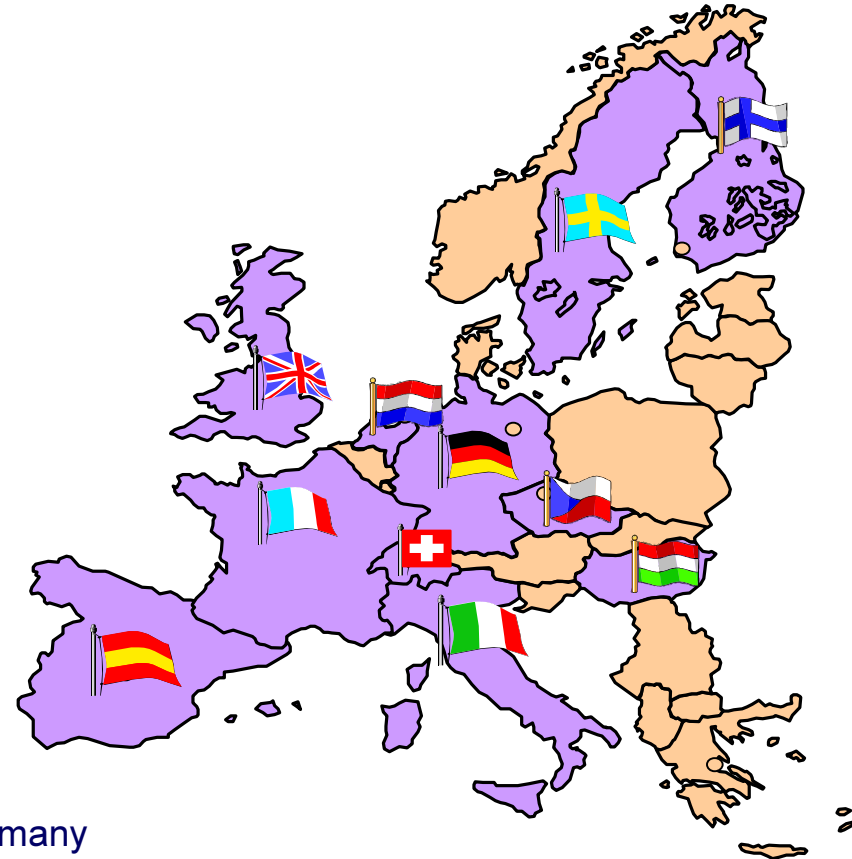
EDG Assistant Partners

Industrial Partners

- Datamat (Italy)
- IBM-UK (UK)
- CS-SI (France)

Research and Academic Institutes

- CESNET (Czech Republic)
- Commissariat à l'énergie atomique (CEA) – France
- Computer and Automation Research Institute, Hungarian Academy of Sciences (MTA SZTAKI)
- Consiglio Nazionale delle Ricerche (Italy)
- Helsinki Institute of Physics – Finland
- Institut de Fisica d'Altes Energies (IFAE) - Spain
- Istituto Trentino di Cultura (IRST) – Italy
- Konrad-Zuse-Zentrum für Informationstechnik Berlin - Germany
- Royal Netherlands Meteorological Institute (KNMI)
- Ruprecht-Karls-Universität Heidelberg - Germany
- Stichting Academisch Rekencentrum Amsterdam (SARA) – Netherlands
- Swedish Research Council - Sweden





EDG overview: Middleware release schedule

- Release schedule
 - testbed 1: 2001
 - testbed 2: 2002
 - testbed 3: 2003
 - Incremental releases between these major dates
- Each release includes
 - feedback on use of previous release by application groups
 - planned improvements/extension by middle-ware groups
- Application groups (HEP, EO, Bio-Info) are using existing software and testbed to explore how they can best exploit grids



EDG overview : current project status

- EDG currently provides a set of middleware services
 - Job & Data Management
 - GRID & Network monitoring
 - Security, Authentication & Authorization tools
 - Fabric Management
- Runs on Linux Red Hat 6.2 + updates platform
 - Site install & config tools and set of common services available
- 5 principle EDG 1.2.0 sites currently belonging to the EDG-Testbed
 - CERN(CH), RAL(UK), NIKHEF(NL), CNAF(I), CC-Lyon(F),
 - being deployed on other EDG testbed sites (~10)
- Intense middleware development continuously going on, concerning:
 - New features for job partitioning and check-pointing, billing and accounting
 - New tools for Data Management and Information Systems.
 - Integration of network monitoring information inside the brokering polices



EDG structure : work packages

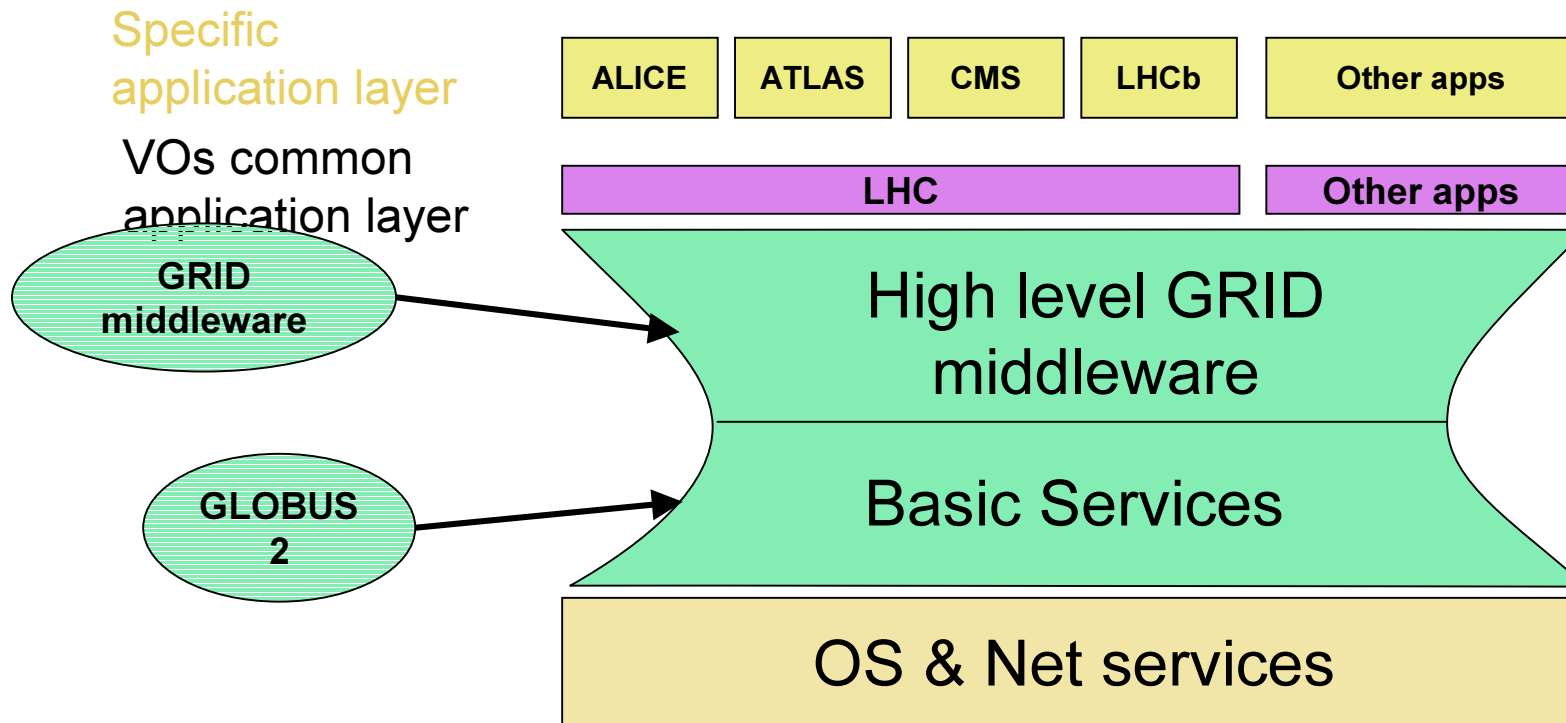
- The EDG collaboration is structured in 12 Work Packages:
 - WP1: Work Load Management System
 - WP2: Data Management
 - WP3: Grid Monitoring / Grid Information Systems
 - WP4: Fabric Management
 - WP5: Storage Element
 - *WP6: Testbed and demonstrators*
 - WP7: Network Monitoring
 - *WP8: High Energy Physics Applications*
 - *WP9: Earth Observation*
 - *WP10: Biology*
 - *WP11: Dissemination*
 - *WP12: Management*
- Applications**



EDG Globus-based middleware architecture

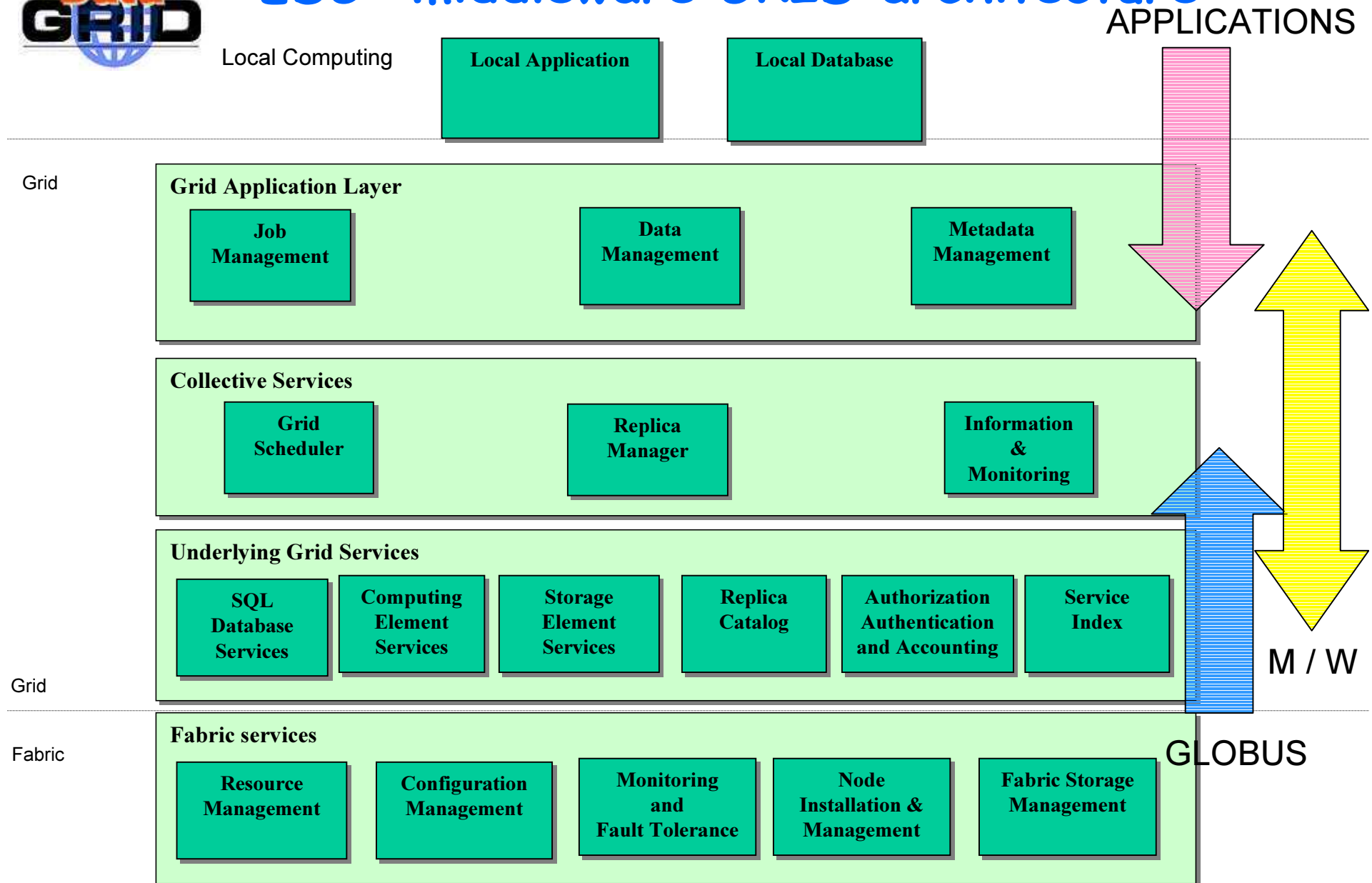
➤ Current EDG architectural functional blocks:

- Basic Services (authentication, authorization, Replica Catalog , secure file transfer, Info Providers) rely on Globus 2
- Higher level EDG middleware (developed within EDG)
- Applications (HEP,BIO,EO)



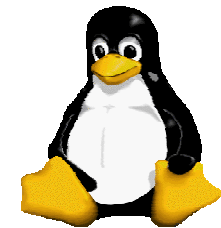
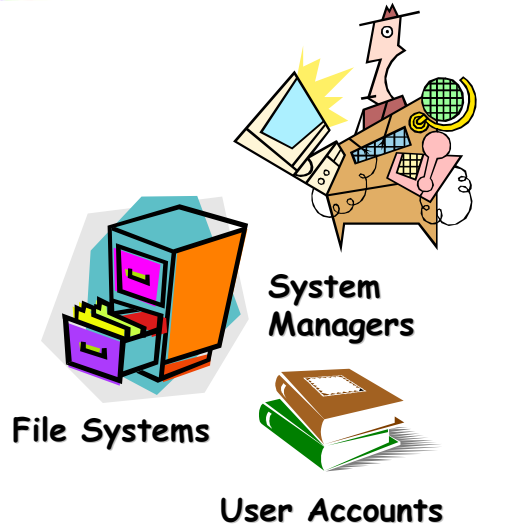


EDG middleware GRID architecture

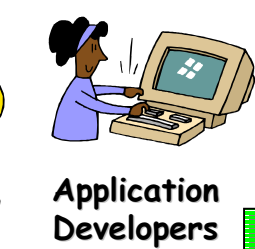




EDG interfaces



Operating Systems



Application Developers



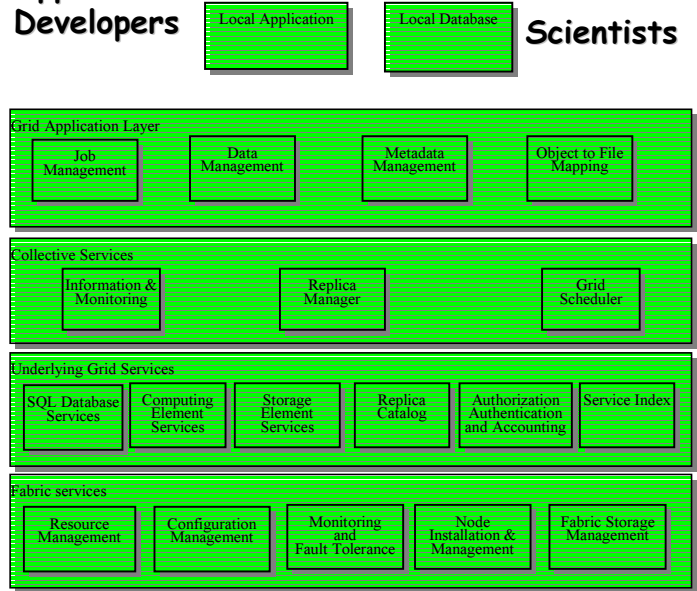
Scientists



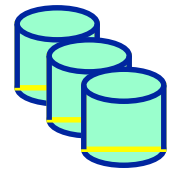
Certificate Authorities



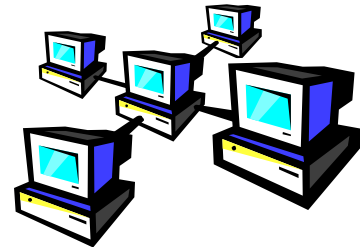
Batch Systems
PBS, LSF, etc.



Mass Storage Systems
HPSS, Castor



Storage Elements



Computing Elements



EDG : reference web sites

- EDG web site
 - <http://www.edg.org>
- Source for all required software :
 - <http://datagrid.in2p3.fr>
- EDG testbed web site
 - <http://marianne.in2p3.fr>
- EDG users guide
 - <http://marianne.in2p3.fr/datagrid/documentation/EDG-Users-Guide.html>
- EDG tutorials web site
 - <http://cern.ch/edg-tutorials>
- EDG production testbed current real time updated set up
 - <http://testbed007.cern.ch/tbstatus-bin/infoindexcern.pl>