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

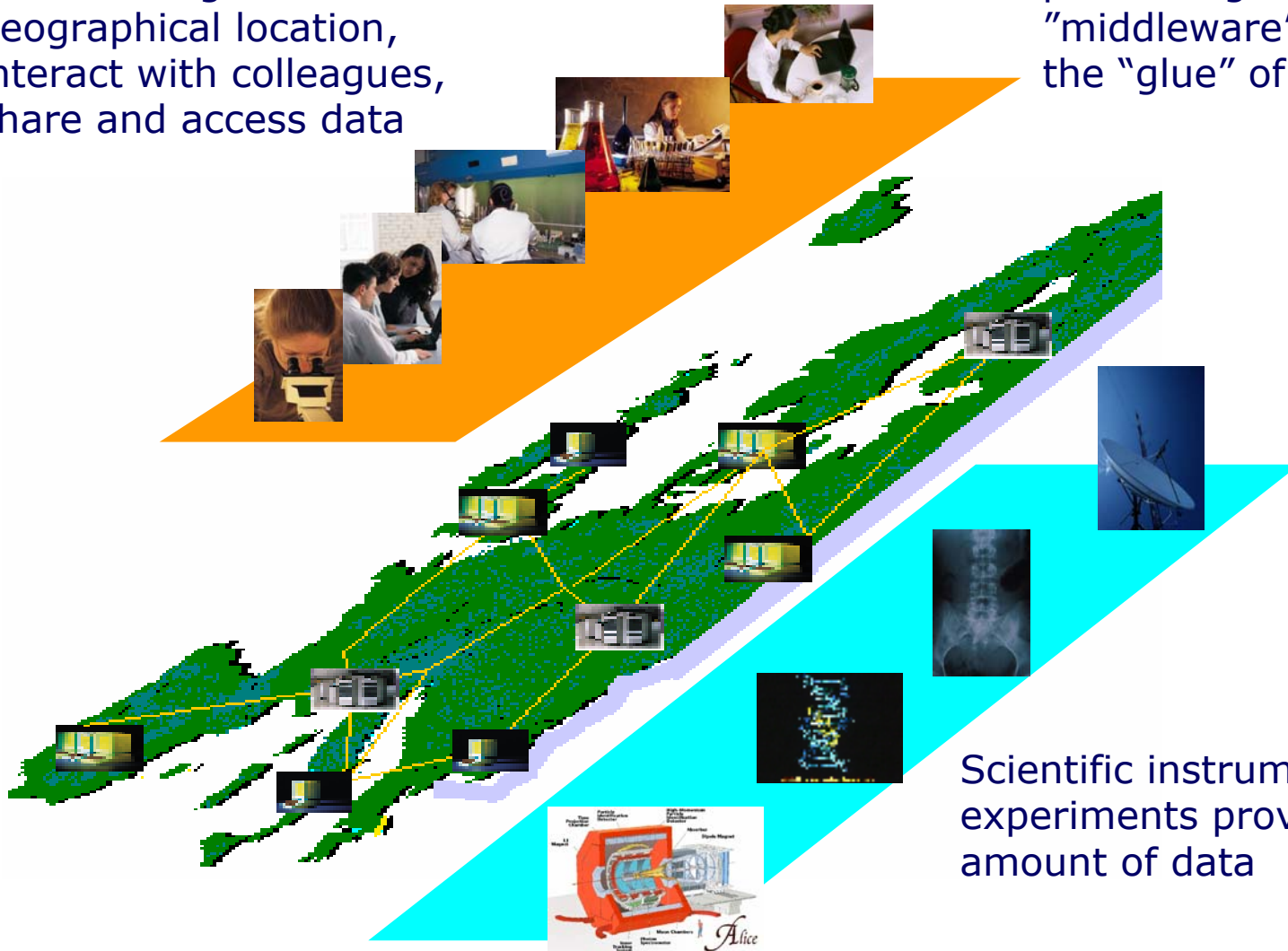**The European DataGrid Project Team**

http://www.eu-datagrid.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

# 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.
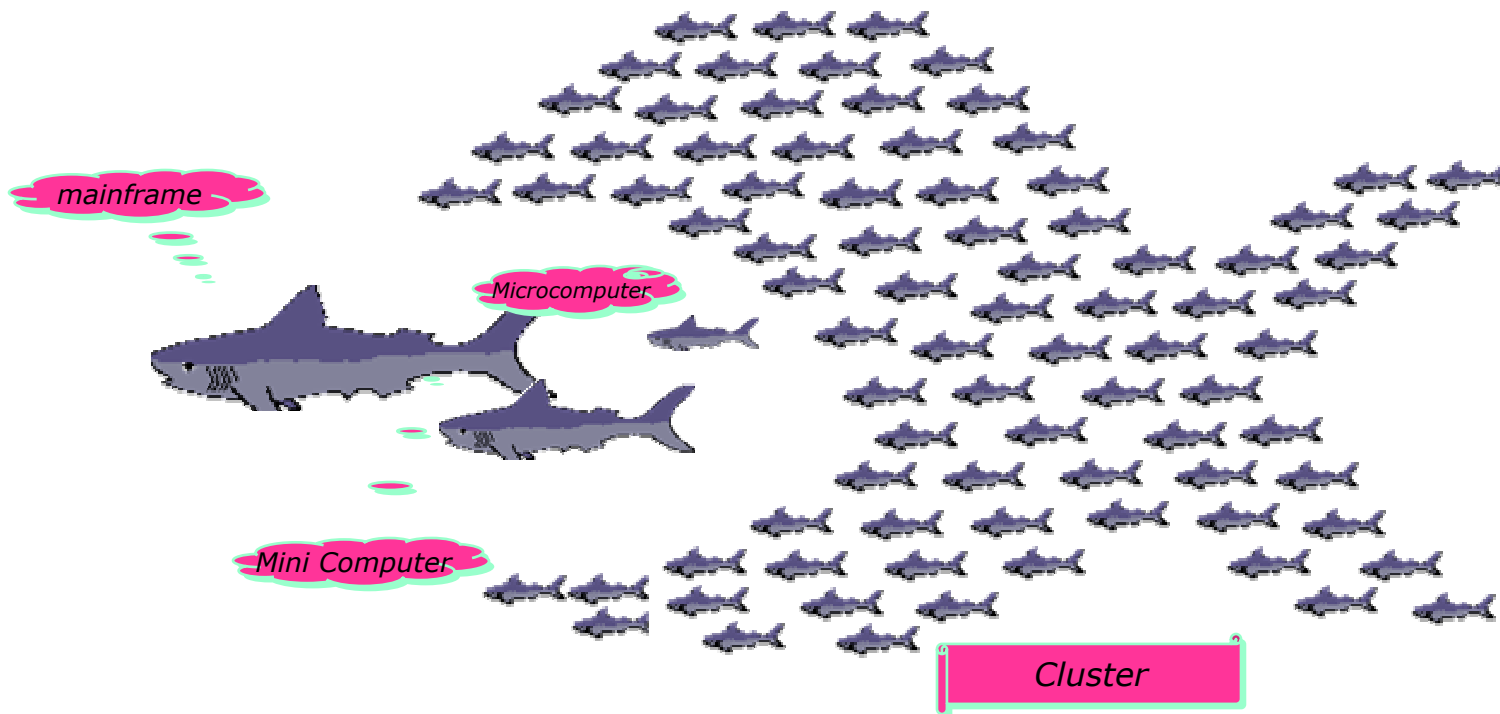
- ◆ Key concept :

  - ▪ **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]

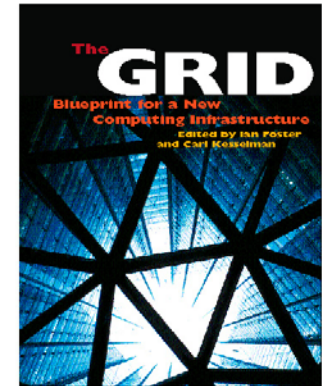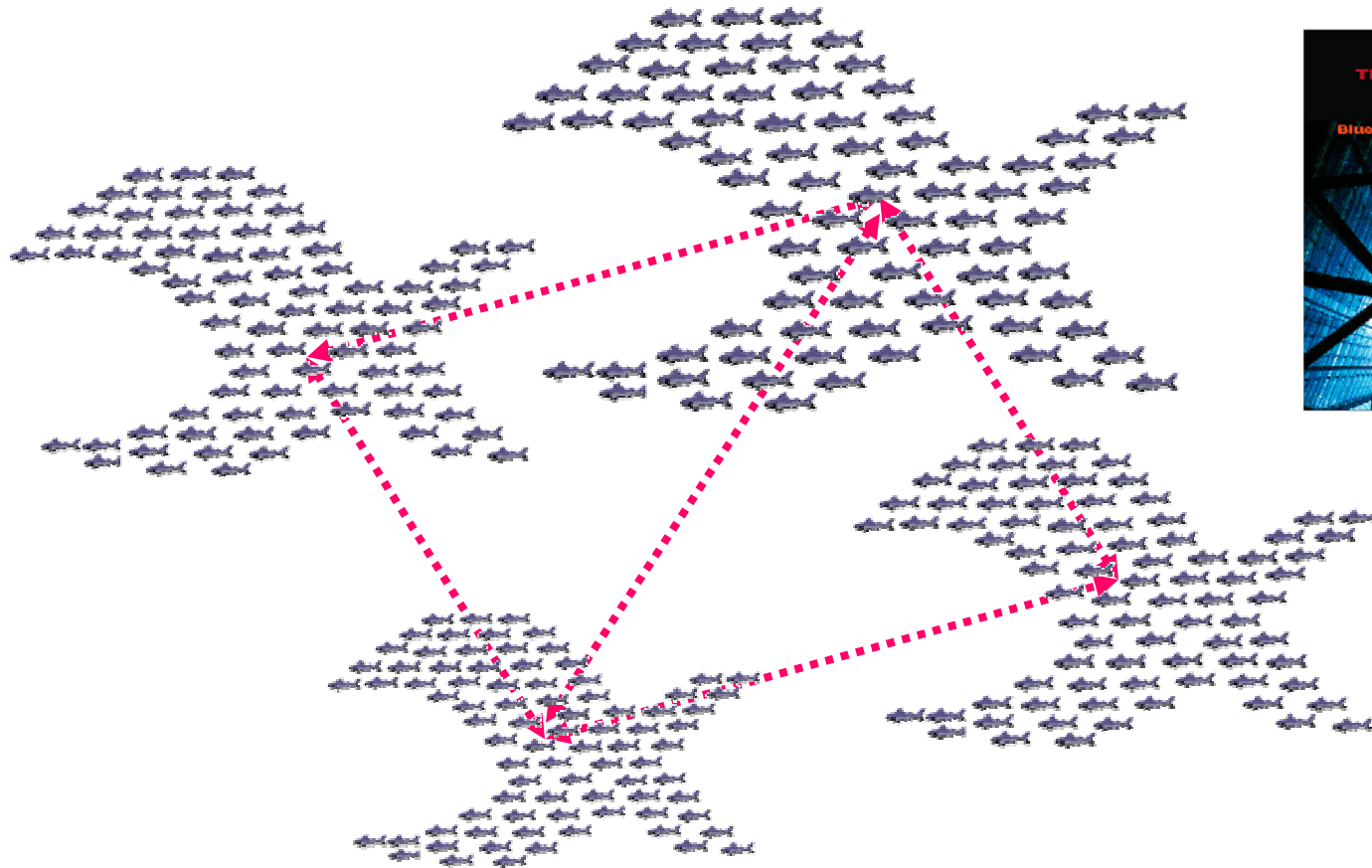# The Grid distributed computing idea 1/2

Once upon a time……..



mainframe

Microcomputer

Mini Computer

Cluster

(by Christophe Jacquet)

# 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 **interactions**
  - User groups or applications want to interact with Grids in different ways

- **Dynamic nature**
  - Resources and users added/removed/changed frequently

# Main Services 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
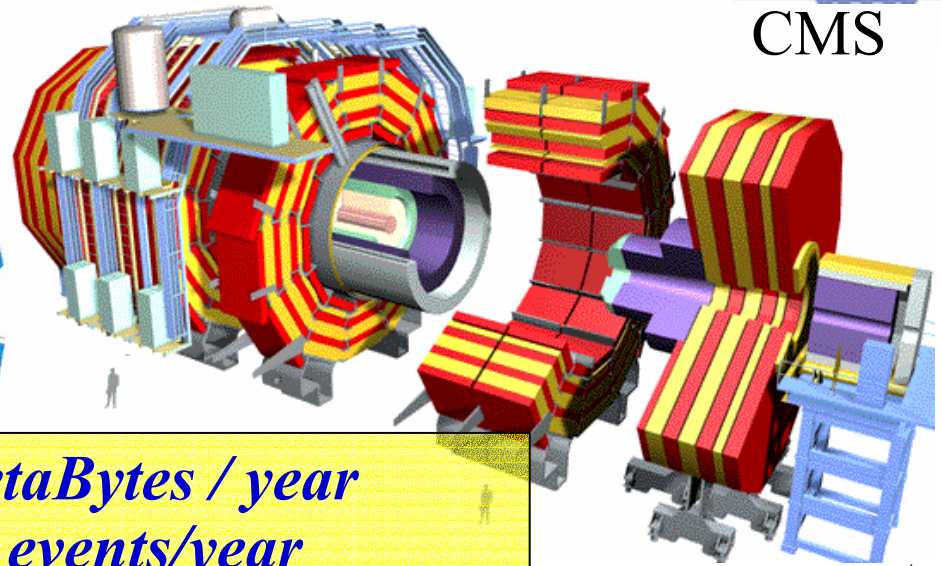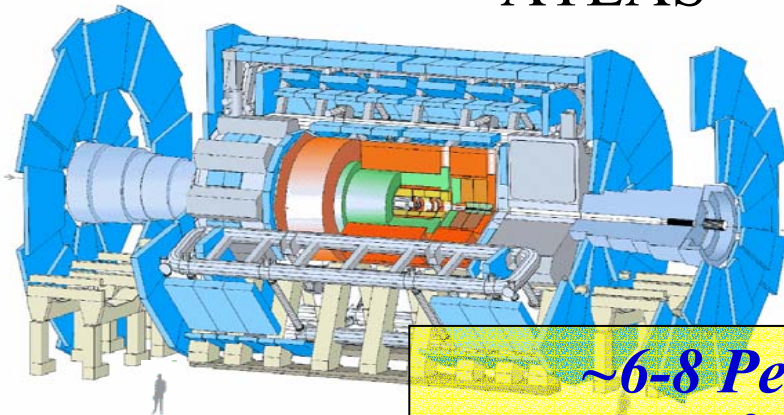
# Why Grids

- ◆ **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

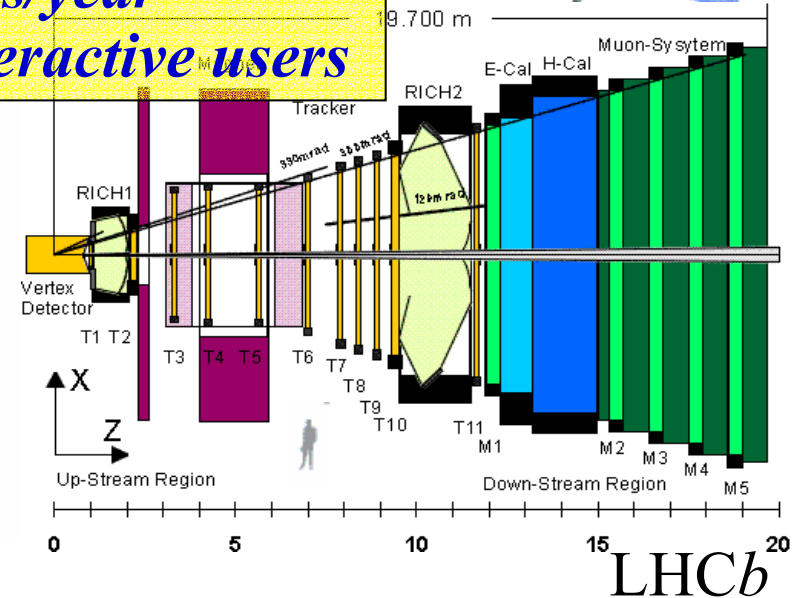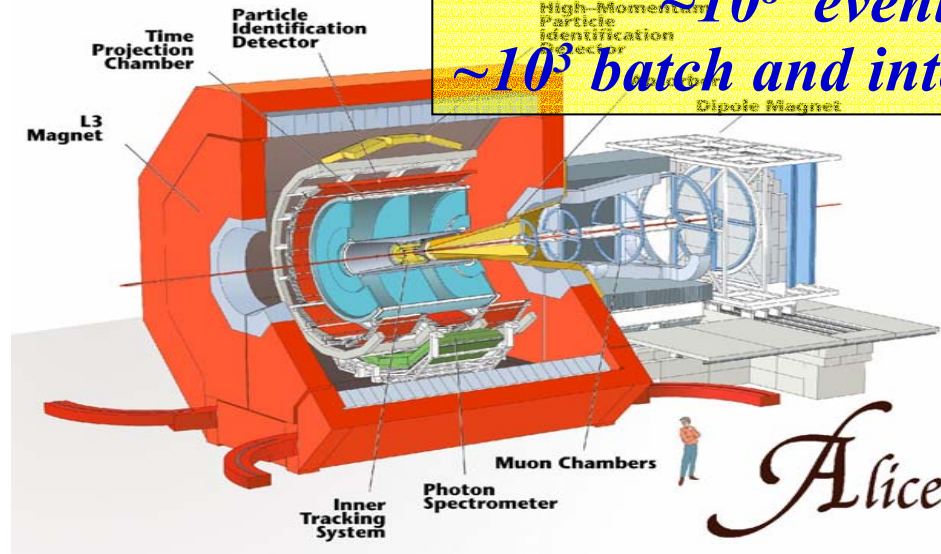  - Astronautics, etc.

# High Energy Physics

*The LHC Detectors*

ATLAS

~6-8 PetaBytes / year
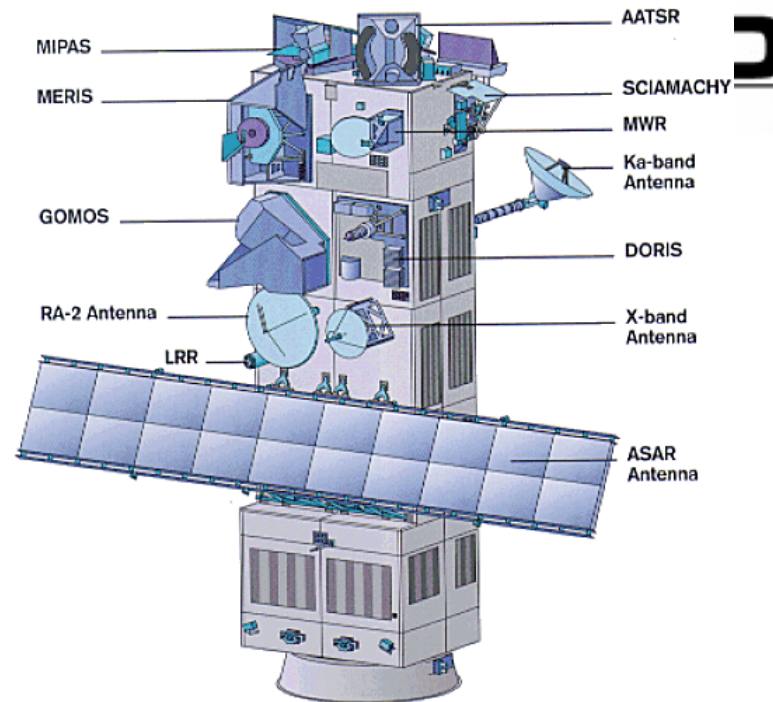~$10^8$ events/year
~$10^3$ batch and interactive users

Alice
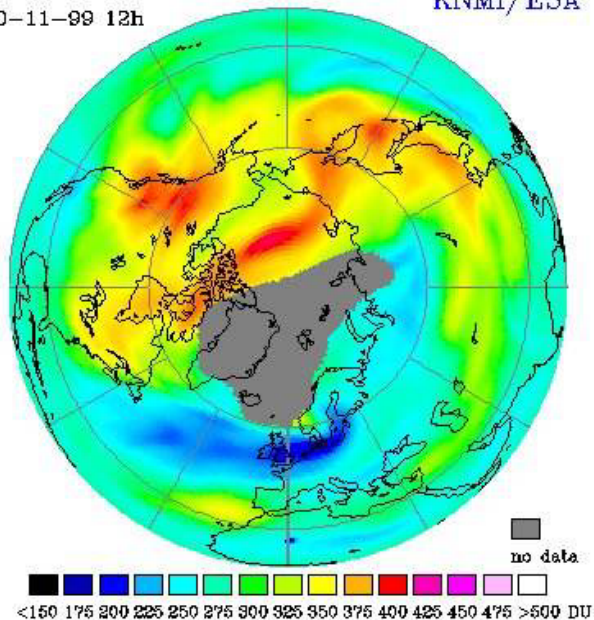
LHC*b*

# 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

<150 175 200 225 250 275 300 325 350 375 400 425 450 475 >500 DU    no data



MIPAS — AATSR
MERIS — SCIAMACHY
— MWR
— Ka-band Antenna
GOMOS —
— DORIS
RA-2 Antenna — X-band Antenna
LRR — 
ASAR Antenna

**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 …)**
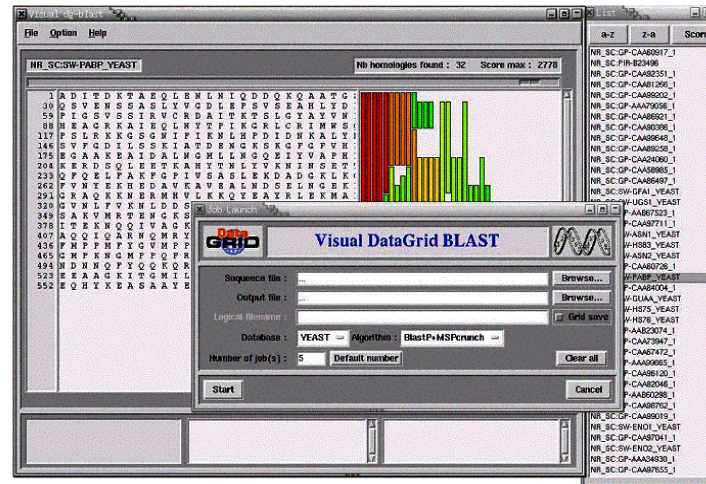
Source: L. Fusco, June 2001

Federico.Carminati , EU review presentation, 1 March 2002

# Biology – BioInformatics

- Bio-informatics
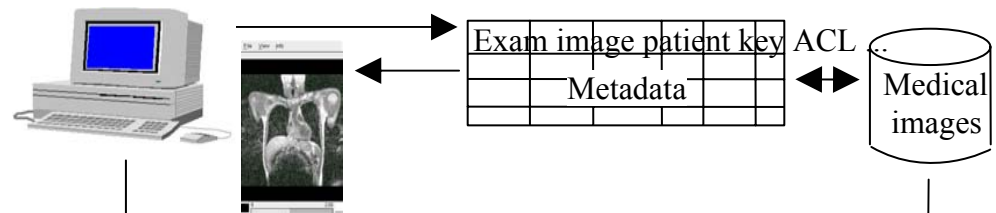  - **Phylogenetics**
  - **Search for primers**
  - **Statistical genetics**
  - **Bio-informatics web portal**
  - **Parasitology**
  - **Data-mining on DNA chips**
  - **Geometrical protein comparison**

- Medical imaging
  - **MR image simulation**
  - **Medical data and metadata management**
  - **Mammographies analysis**
  - **Simulation platform for PET/SPECT**

| | |
|---|---|
| 🟩 | **Applications deployed** |
| 🟧 | **Applications tested on EDG** |
| 🟥 | **Applications under preparation** |



1. Query the medical image database and retrieve a patient image

Exam image patient key ACL

Metadata

Medical images

2. Compute similarity measures over the database images

Submit 1 job per image

3. Retrieve most similar cases

Similar images          Low score images

# Major existing Grid projects (1/2)

◆ **Europe-based projects:**

- **European DataGrid (EDG) : 2001-2003**      **www.edg.org**

- **LHC Computing GRID (LCG): 2002-2008 -....**     **cern.ch/lcg**

- **CrossGrid**                   **: 2002-2005 www.crossgrid.org**

- **DataTAG**                    **: 2002-2003 www.datatag.org**

- **GridLab**                    **: 2002-2004**      **www.gridlab.org**

- **EGEE**                     **: 2004-2007 ?**      **www.cern.ch/egee**

**European National Projects:**

- **INFNGRID, UK-GridPP, NorduGrid(Nordic test bed for wide area computing )…**

# Major existing Grid projects (2/2)

- **US projects:**

  - GriPhyN   HEP    www.griphyn.org

  - PPDG       HEP    www.ppdg.net

  - iVDGL  ( joint GriPhyN, PPDG)  www.ivdgl.or

  - TERAGRID (NSF)        www.teragrid.org
    - IBM, Intel Qwest ,Myricom, Sun Microsystems, Oracle.

  - National Middleware Initiative (NSF NMI)    www.nsf-middleware.org

  - ESG        www.earthsystemgrid.org

  - NEESgrid   virtual lab earthquake engineering   www.neesgrid.org

  - BIRN biomedical informatics research network  birn.ncrr.nih.gov/birn/

- **Asia-based projects**:

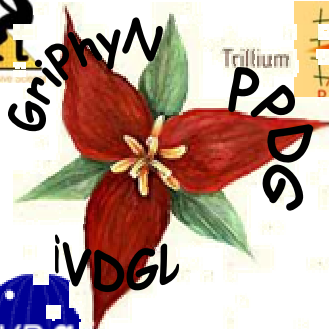  - ApGRID         www.apgrid.org

  - TWGRID         www.twgrid.org

  - Many Grid projects in : Korea, Japan, China,    Australia

# Major US & European Grid Projects,
## many with strong HEP participation



**The Virtual Data Toolkit (VDT)**

GriPhyN

PPDG

iVDGL

*US projects*

DataTAG

Many national, regional Grid projects ·
GridPP(UK), INFN-grid(I), NorduGrid, Dutch Grid, …

**The DataGrid Toolkit**

*European projects*

# The European Data Grid Project

- ◆ 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

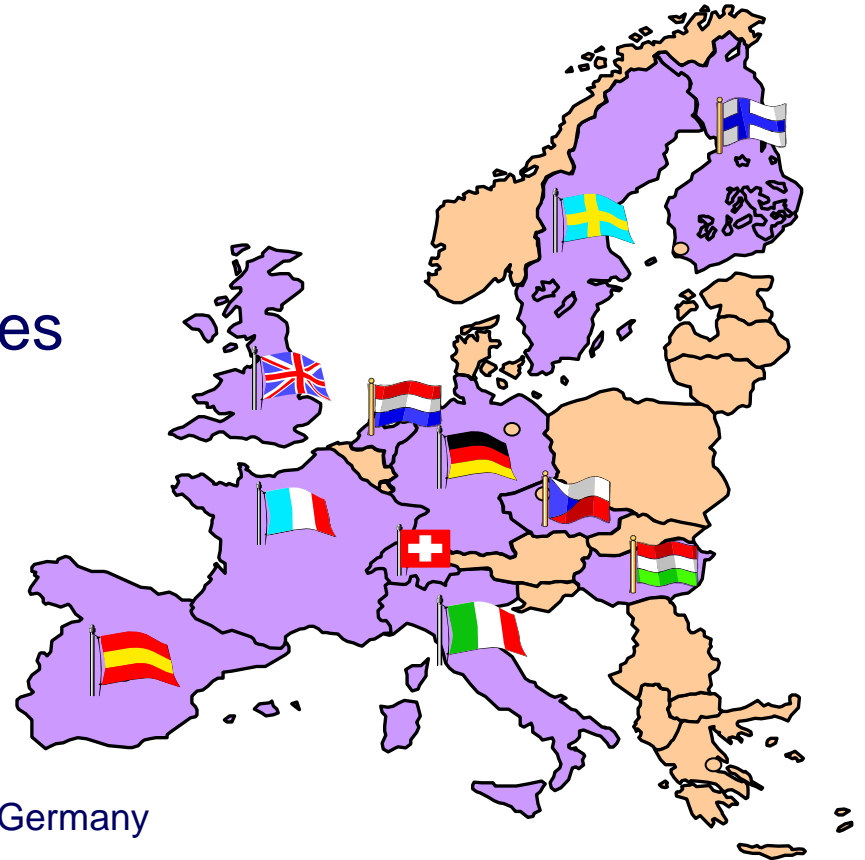- ➢ PPARC - UK

# 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:** late 2001
  - **testbed 2:** early 2003
  - **testbed 3:** end 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

# Current Project Status

- EDG currently provides a set of middleware services
    - Job & Data Management
    - Grid & Network monitoring
    - Security, Authentication & Authorization tools
    - Fabric Management

- EDG release 2.0 currently deployed to the EDG-Testbeds
    - GNU/Linux RedHat 7.3 on Intel PCs ~15 sites in application testbed actively used by application groups
        - Core sites CERN(CH), RAL(UK), NIKHEF(NL), CNAF(I), CC-Lyon(F)
    - EDG sw also deployed at total of ~40 sites via CrossGrid, DataTAG and national grid projects

- Final release 2.1 will be out soon

- Many applications ported to EDG testbeds and actively being used

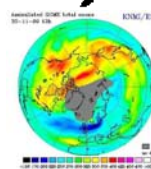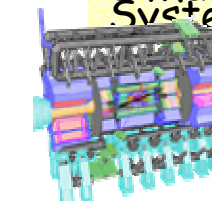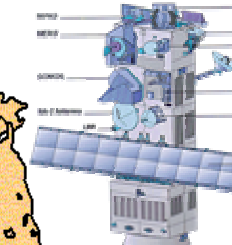- Intense middleware development continuously going-on

# DataGrid in Numbers

**People**

>350 registered users

12 Virtual Organisations

16 Certificate Authorities

>500 people trained

278 man-years of effort

100 years funded

**Software**

50 use cases

18 software releases

>300K lines of code

**Testbeds**

>15 regular sites

>10'000s jobs submitted

>1000 CPUs

>5 TeraBytes disk

3 Mass Storage Systems

**Scientific applications**

5 Earth Obs institutes

9 bio-informatics apps

6 HEP experiments

# 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
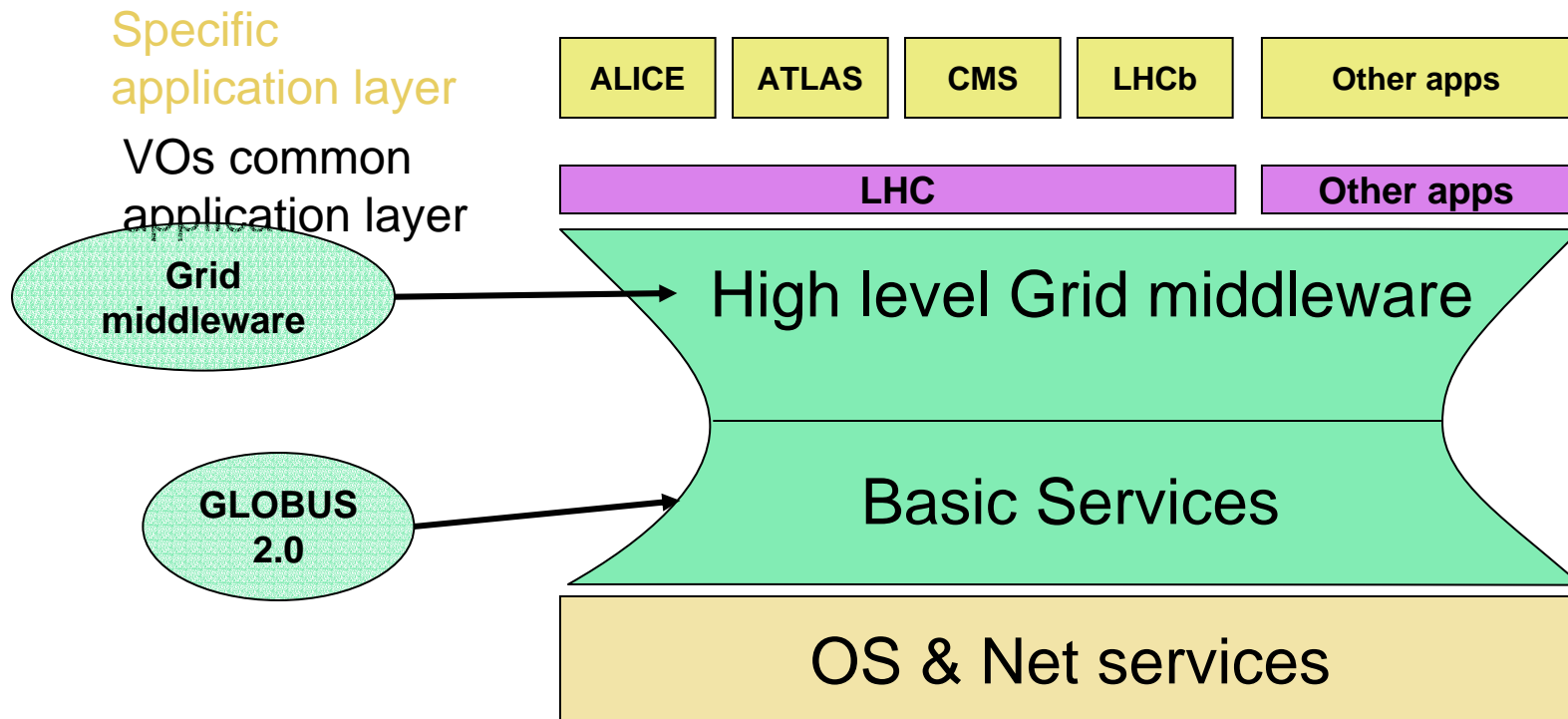
- WP11: Dissemination
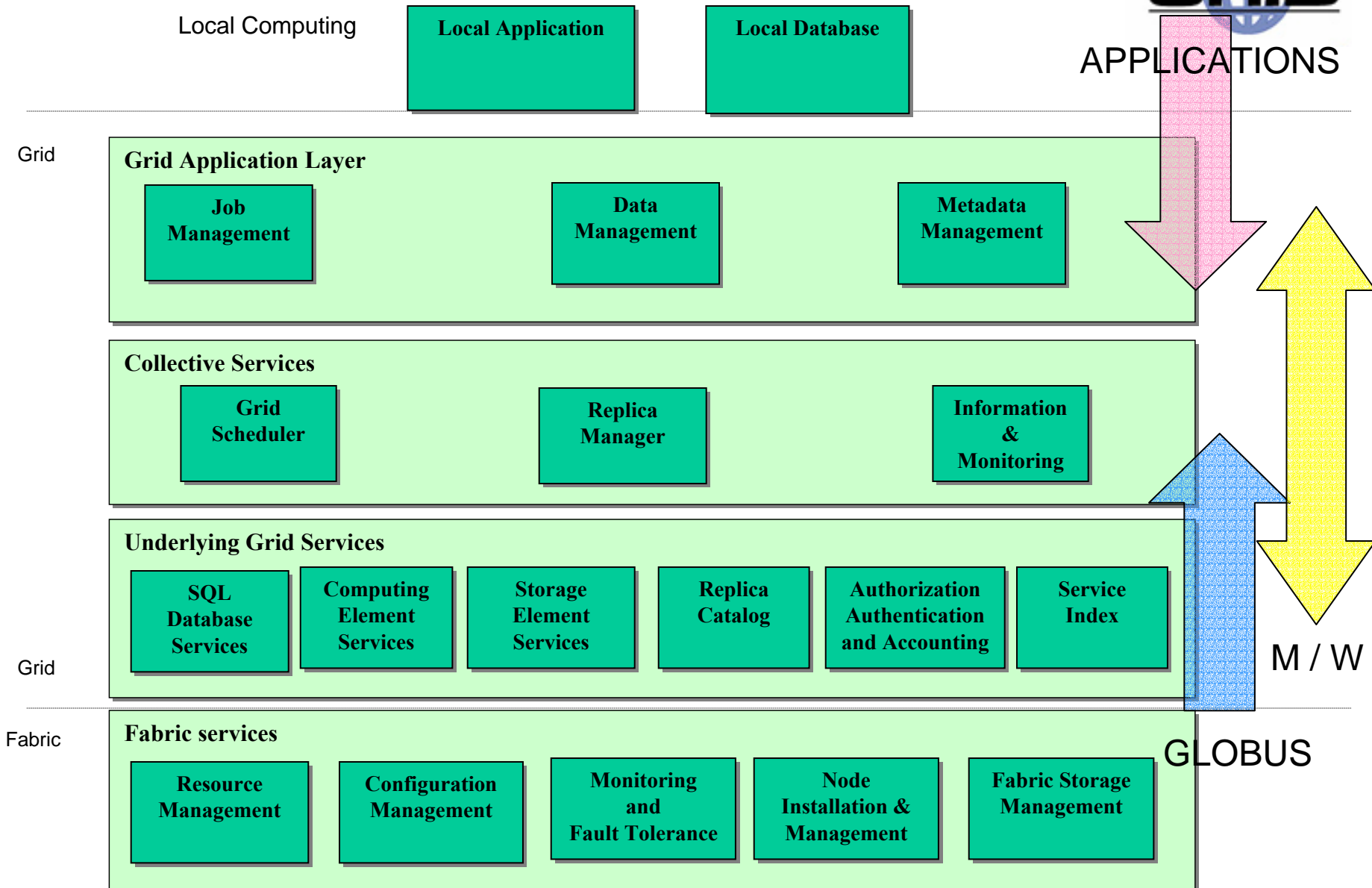
- WP12: Management

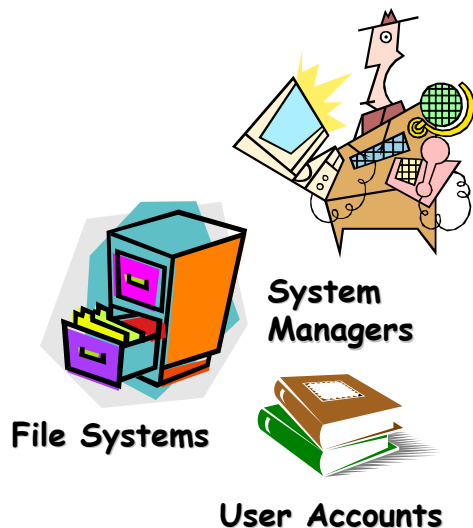**Applications**

# EDG Globus-based middleware architecture

➤ Current EDG architectural functional blocks:

  ▪ **Basic Services** (authentication, authorization, Replica Catalog , secure file transfer, Info Providers) rely on Globus 2.0

  ▪ **Higher level EDG middleware.**(developed within EDG)

  ▪ **Applications** (HEP,BIO,EO)

Specific application layer

| ALICE | ATLAS | CMS | LHCb | Other apps |

VOs common application layer

| LHC | Other apps |

Grid middleware → High level Grid middleware

GLOBUS 2.0 → Basic Services

OS & Net services

# EDG middleware Grid architecture

**Local Computing**

| Local Application | Local Database |
|---|---|

**APPLICATIONS**

**Grid**

## Grid Application Layer

| Job Management | Data Management | Metadata Management |
|---|---|---|

## Collective Services

| Grid Scheduler | Replica Manager | Information & Monitoring |
|---|---|---|

## Underlying Grid Services

| SQL Database Services | Computing Element Services | Storage Element Services | Replica Catalog | Authorization Authentication and Accounting | Service Index |
|---|---|---|---|---|---|

**Grid**

**M / W**

**Fabric**

## Fabric services

| Resource Management | Configuration Management | Monitoring and Fault Tolerance | Node Installation & Management | Fabric Storage Management |
|---|---|---|---|---|

**GLOBUS**

# EDG Interfaces



Application Developers

Local Application

Local Database

Scientists

Certificate Authorities

System Managers

File Systems

User Accounts

**Grid Application Layer**
| Job Management | Data Management | Metadata Management | Object to File Mapping |

**Collective Services**
| Information & Monitoring | Replica Manager | Grid Scheduler |

**Underlying Grid Services**
| SQL Database Services | Computing Element Services | Storage Element Services | Replica Catalog | Authorization Authentication and Accounting | Service Index |

**Fabric services**
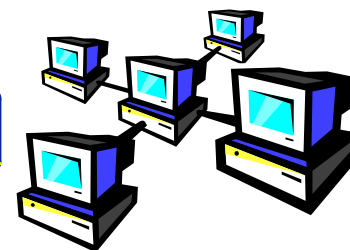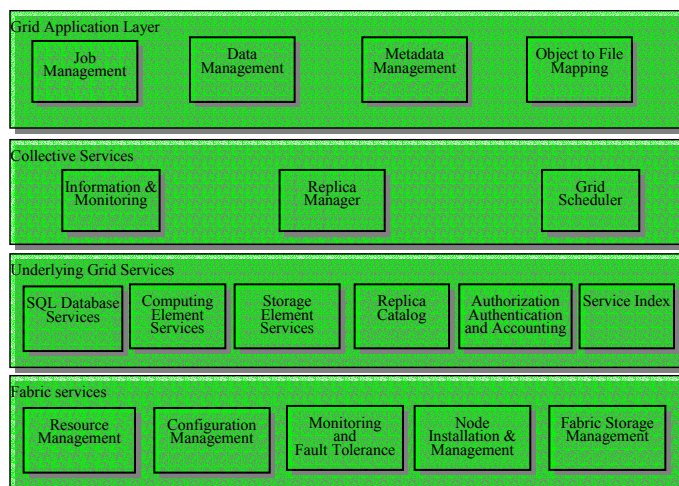| Resource Management | Configuration Management | Monitoring and Fault Tolerance | Node Installation & Management | Fabric Storage Management |

Operating Systems

Mass Storage Systems
HPSS, Castor

Storage Elements

Computing Elements

Batch Systems
PBS, LSF, etc.

the globus project™
www.globus.org

Condor
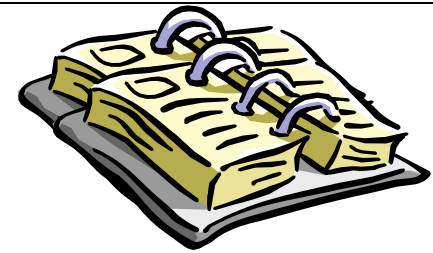High Throughput Computing

SECURITY
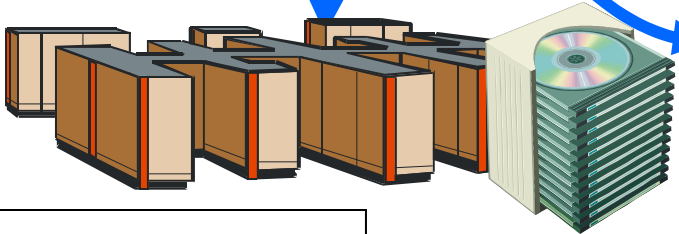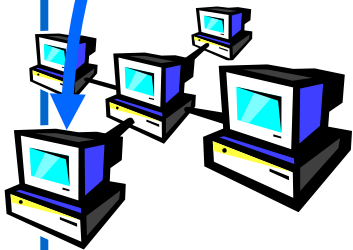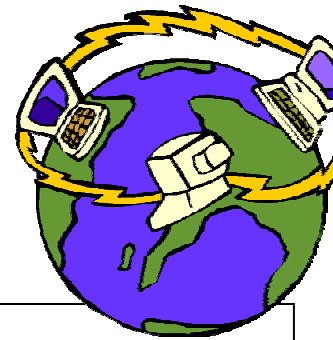
# EDG Tutorial Overview



Workload Management Services

Data Management Services

Networking

Information Service

Fabric Management

# 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

- ◆ Dissemination Testbed (GriDis)
  - **http://web.datagrid.cnr.it/GriDis/GriDisWP1.html**

- ◆ EDG users guide
  - http://marianne.in2p3.fr/datagrid/documentation/EDG-Users-Guide.html

- ◆ EDG tutorials web site
  - http://cern.ch/edg-tutorials