# Overview of the EU Data Grid Project

**The European
DataGrid Project Team**

http://www.eu-datagrid.org

# The EU DataGrid (EDG) Project

- ◆ 9.8 M Euros EU (IST) funding over 3 years

- ◆ Three year phased developments & demos (2001-2003)

- ◆ 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

- ◆ Total of 21 partners
  - ▪ Research and Academic institutes as well as industrial companies

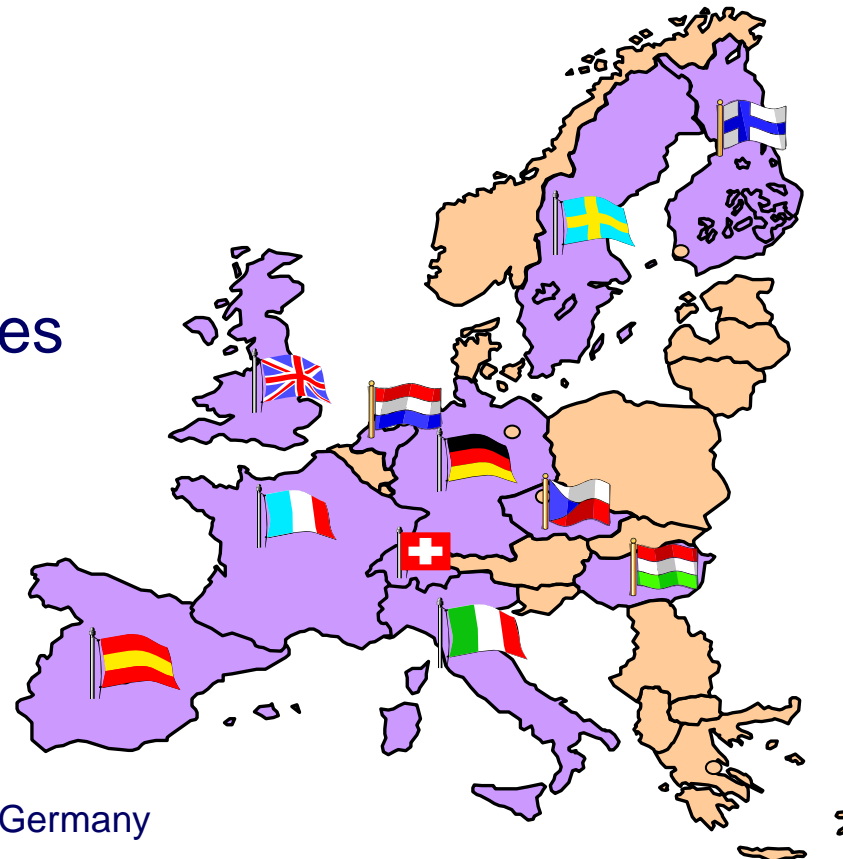- ◆ Main partners:

# 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 structure : work packages

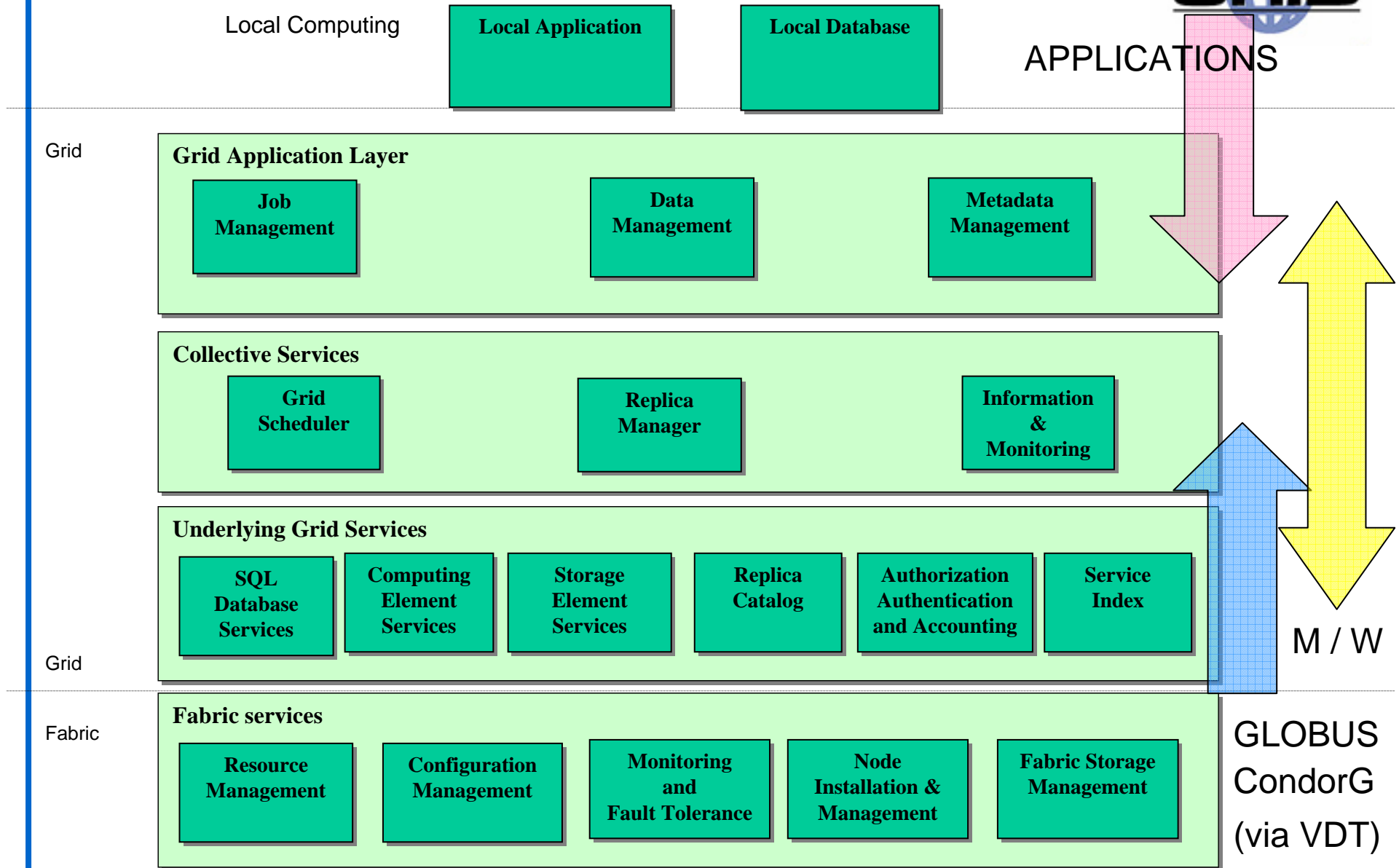> The EDG collaboration is structured in 12 Work Packages:

- 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

> **Applications**

- WP11: Dissemination

- WP12: Management

# Project Schedule

- Project started on 1/Jan/2001

- Testbed 0 (early 2001)
  - International test bed 0 infrastructure deployed
    - Globus 1 only - no EDG middleware

- Testbed 1 ( 2002 )
  - First release of EU DataGrid software to defined users within the project:
    - HEP experiments (WP 8), Earth Observation (WP 9), Biomedical applications (WP 10)

- Testbed 2 (End 2002)
  - Builds on Testbed 1 to extend facilities of DataGrid
  - Focus on stability

- EDG very successfully passed its 2nd annual EU review on February 4-5 2003

- Testbed 3 (2003)
  - Advanced functionality; currently being deployed.

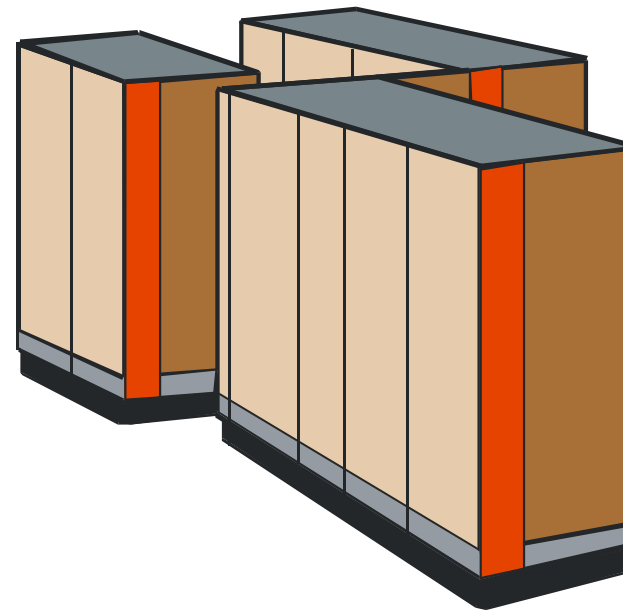- Project stops on 31/Dec/2003

# EDG Middleware 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

### Fabric services

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

Fabric

M / W

GLOBUS
CondorG
(via VDT)

# The EDG WMS (WP1)

- The user interacts with Grid via a **Workload Management System** (WMS)

- The Goal of WMS is the **distributed scheduling and resource management in a Grid environment**.

- What does it allow Grid users to do?

  - To submit their jobs

  - To execute them on the "best resources"
    - The WMS tries to optimize the usage of resources

  - To get information about their status
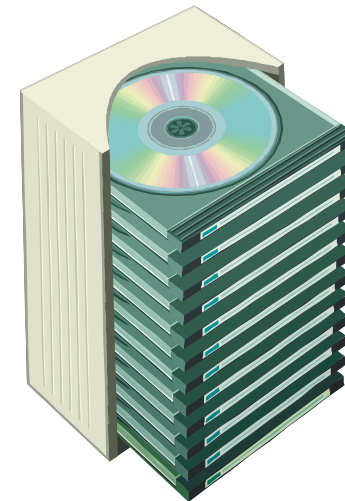
  - To retrieve their output

# Computing Element (WP4)

◆ Is a Grid Job Queue

- Publishes information about itself

- Checks the job is permitted

- Sends it to an an appropriate internal queue

# SRM: Storage Resource Manager (WP5)

◆ SRM subset implementation

   ▪ A defacto international standard for Storage Resource Management

◆ Web service uses Java AXIS and EDG security

◆ Supports multiple VOs

◆ Functions

   ▪ Writing a file

   ▪ Reading a file

# Replica Manager (WP2)

- ◆ High level data management on the Grid
  - ▪ Location of data
  - ▪ Replication of data
  - ▪ Efficient access to data

- ◆ Hides the SRM

- ◆ Coordinates use of
  - ▪ Replica Location Service
  - ▪ Replica Metadata Catalog
  - ▪ Replica Optimization Service

# Information & Monitoring (WP3) R-GMA

◆ Relational implementation of GMA from GGF

◆ Makes use of GLUE schema

◆ Interoperable with MDS

◆ Deals with information on

- The Grid itself
  - Resources and Services (for which the Globus MDS is a common solution)
  - Job status information
- Grid applications
  - This is information published by user jobs.

# TestBed Integration (WP6)

- Exact definition of RPM lists (components) for the various testbed machine profiles (CE, RB, UI,, WN, IC etc.) – check dependencies

- Perform preliminary centrally (CERN) managed tests on EDG m/w before green light for spread EDG testbed sites deployment

- Provide, update end user documentation for installers/site managers, developers  and end users

- Define EDG release policies, coordinate the integration team staff with the various WorkPackage managers – keep high inter-coordination.

- Set up the Authorization Working Group to manage authorization policies on the testbed



EDG production testbed  EDG Release EDG 1.4.7
CERN Computing Element site and corresponding spread GRID services

# Grid aspects covered by EDG

| | | | |
|---|---|---|---|
| **VOMS** | Provides certificate with VOs, groups and roles | **RGMA: Information & Monitoring** | Provides info on resource utilization & performance |
| **User Interface** | Submit & monitor jobs, retrieve output | **Grid Fabric Management** | Configure, installs & maintains grid sw packages and environ. |
| **Workload Management System** | Manages submission of jobs to Res. Broker, obtains information and retrieves output | **Network performance** | Provides efficient network transport, bandwidth monitoring |
| **Computing Element** | Gatekeeper to a grid computing resource | **Testbed admin.** | Certificate auth.,user reg., usage policy etc. |
| **Storage Resource Manager** | Grid-aware storage area | **Applications** | HEP, EO, Biology |
| **Replica Manager** | Replicates and locates data | | |

# EDG Interfaces



**Application Developers**

**Scientists**

**System Managers**

**File Systems**

**User Accounts**

**Certificate Authorities**

Local Application

Local Database

| 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 |

the globus project®
www.globus.org

**Condor**
*High Throughput Computing*

**SECURITY**

**Operating Systems**

**Mass Storage Systems HPSS, Castor**

**Storage Elements**

**Computing Elements**

**Batch Systems PBS, LSF, etc.**

# DataGrid in Numbers (as of Feb. 2003)

**People**

>350 registered users

12 Virtual Organisations

16 Certificate Authorities

>200 people trained

278 man-years of effort

   100 years funded

**Testbeds**

>15 regular sites

>40 sites using EDG sw

>10'000s jobs submitted

>1000 CPUs

>15 TeraBytes disk

3 Mass Storage Systems

**Software**

50 use cases

18 software releases

Current release 1.4

>300K lines of code

**Scientific applications**

5 Earth Obs institutes

9 bio-informatics apps

6 HEP experiments

# DataGrid Scientific Applications

*Developing grid middleware to enable large-scale usage by scientific applications*



## Bio-informatics

- Data mining on genomic databases (exponential growth)
- Indexing of medical databases (Tb/hospital/year)



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

no data

<150 175 200 225 250 275 300 325 350 375 400 425 450 475 >500 DU

## Earth Observation

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



## Particle Physics

- Simulate and reconstruct complex physics phenomena millions of times
- LHC experiments will generate 6-8 PetaBytes/year

# Application Usage of Release 1.4

**EDG 1.4 evaluated for review in Feb. 2003**
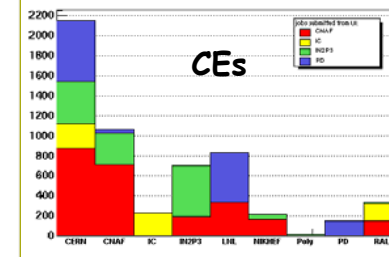
**CPU Usage**



## Positive Signs:

- Large increase in users.

- Many sites interested in joining.

- Pushing real jobs through system.



**Disk Usage (CERN)**

**TOTAL: >1.5 TB**

**HEP Simulation**

**Disk Usage**



**SEs**

**CEs**

# Release 2.0

- ◆ Major new developments in all middleware areas

- ◆ Addre **Being deployed on Application TB now**
  - ▪ WMS **Tutorial will use 1.4 the first day**
    - Job submission
  - ▪ Rep
    - Data management
  - ▪ Data
  - ▪ Info **2.0 will be used second day**
    scal
    - new job submission features
  - ▪ Unif
    - new data mgmt features
  - ▪ Fabr
    - R-GMA

- ◆ Provi

- ◆ Upgrade underlying software

# Related Grid Projects



**Through links with sister projects, there is the potential for a truly global scientific applications grid**
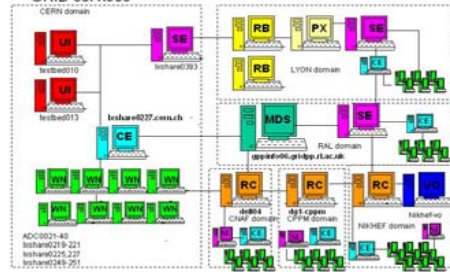**Demonstrated at IST2002 and SC2002 in November**

# Tutorial Testbed

- EDG dissemination testbed (GriDis)

- Resources from EDG and Crossgrid

# EDG Tutorial Roadmap

**Security**

**Testbed**

EDG production testbed  EDG Release EDG 1.4.7
CERN Computing Element site and corresponding spread
GRID services

**Workload Management**

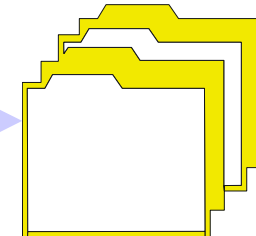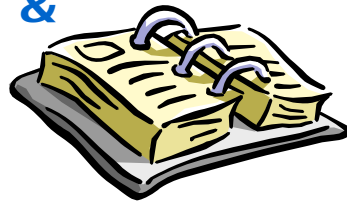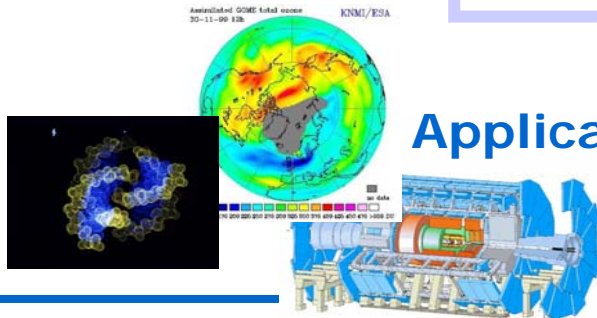**Data Management**

**Information & Monitoring**

**Applications**

**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