# The EGEE project: building an international production grid infrastructure

## "A Biomed view"

- **EGEE -** *what is it and why is it needed?*
- **Middleware –** *current and future*
- **Operations –** *providing a stable service*
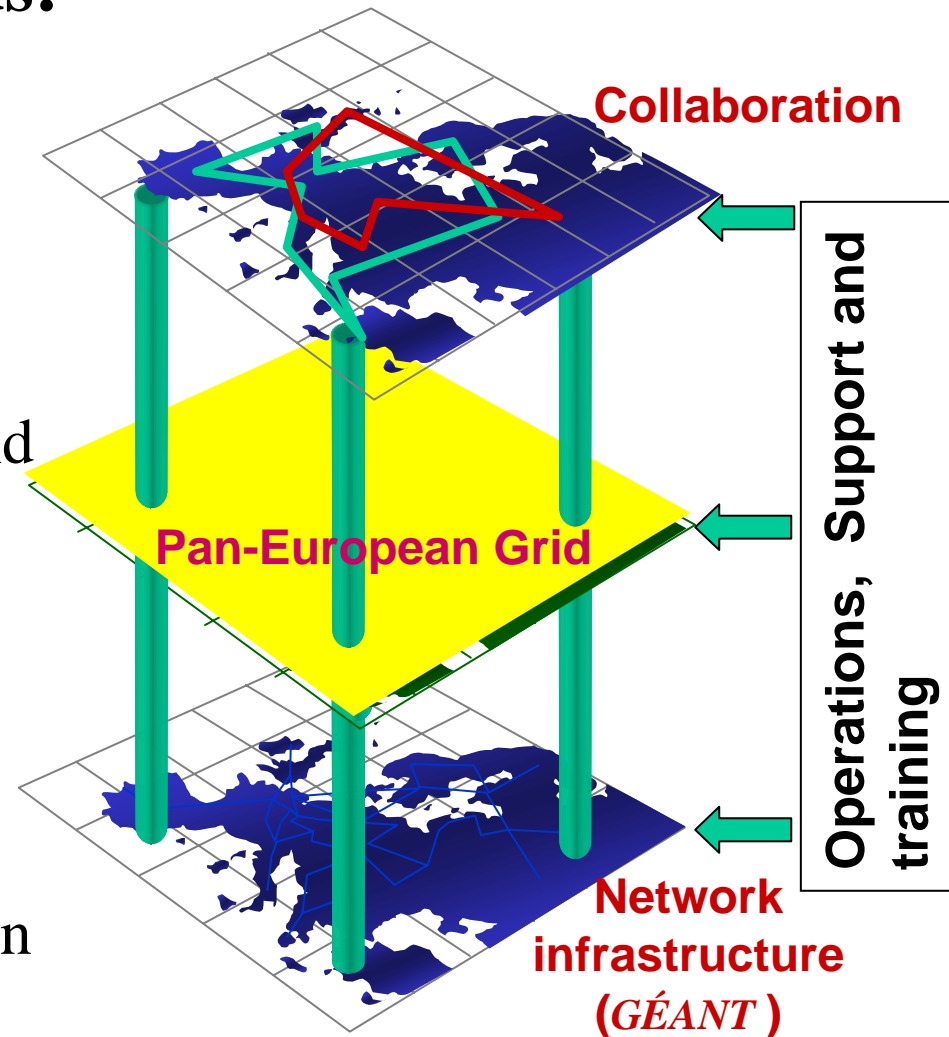- **Networking –** *enabling collaboration*
- **Summary**

*The material for this talk has been contributed by many colleagues in the EGEE & LCG projects.*

*It is heavily based on Bob Jones' talk at UK AHM 2004.*

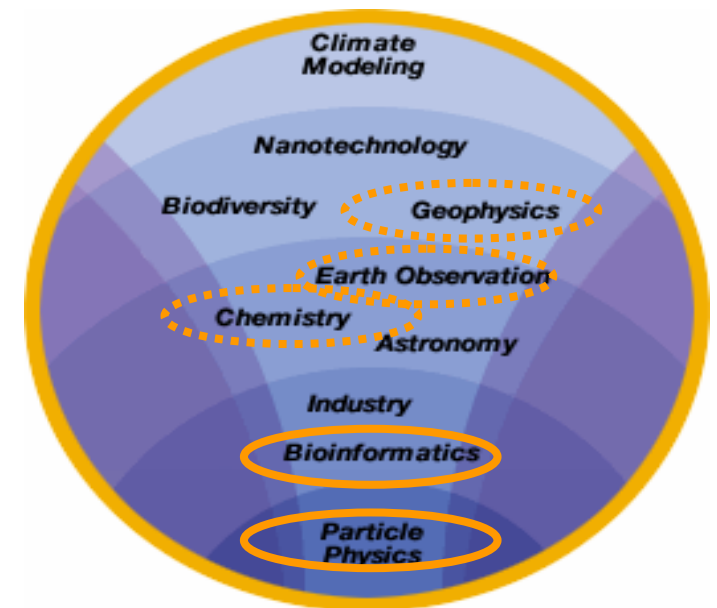# The next generation of grids: EGEE Enabling Grids for E-science

**Build a large-scale production grid service to:**

- Underpin European science and technology

- Link with and build on national, regional and international initiatives

- Foster international cooperation both in the creation and the use of the e-infrastructure



**Collaboration**

**Pan-European Grid**

**Operations, Support and training**

**Network infrastructure (*GÉANT*)**

Information Society

# In 2 years EGEE will:

- **Establish production quality sustained Grid services**
  - 3000 users from at least 5 disciplines
  - over 8,000 CPU's, 50 sites
  - over 5 Petabytes ($10^{15}$) storage

- Demonstrate a viable general process to **bring other scientific communities on board**



- **Propose a second phase** in mid 2005 to take over EGEE in early 2006

- **Frank finds that EGEE has a focus on Biomed applications**

- **It provides support and already has a Biomed VO**
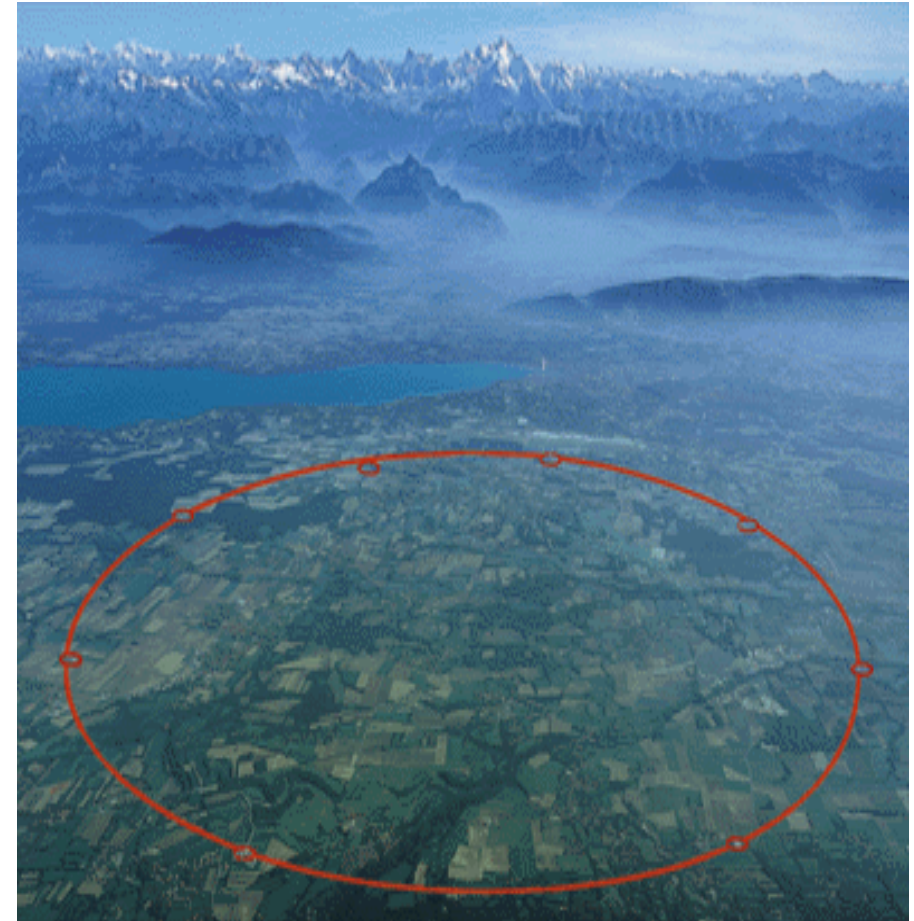
Information Society

EGEE builds on the work of LCG to establish a grid operations service

- **LCG (LHC Computing Grid) - Building and operating the LHC Grid**
- A collaboration between:
  - The physicists and computing specialists from the LHC experiment
  - The projects in Europe and the US that have been developing Grid middleware
  - The regional and national computing centres that provide resources for LHC
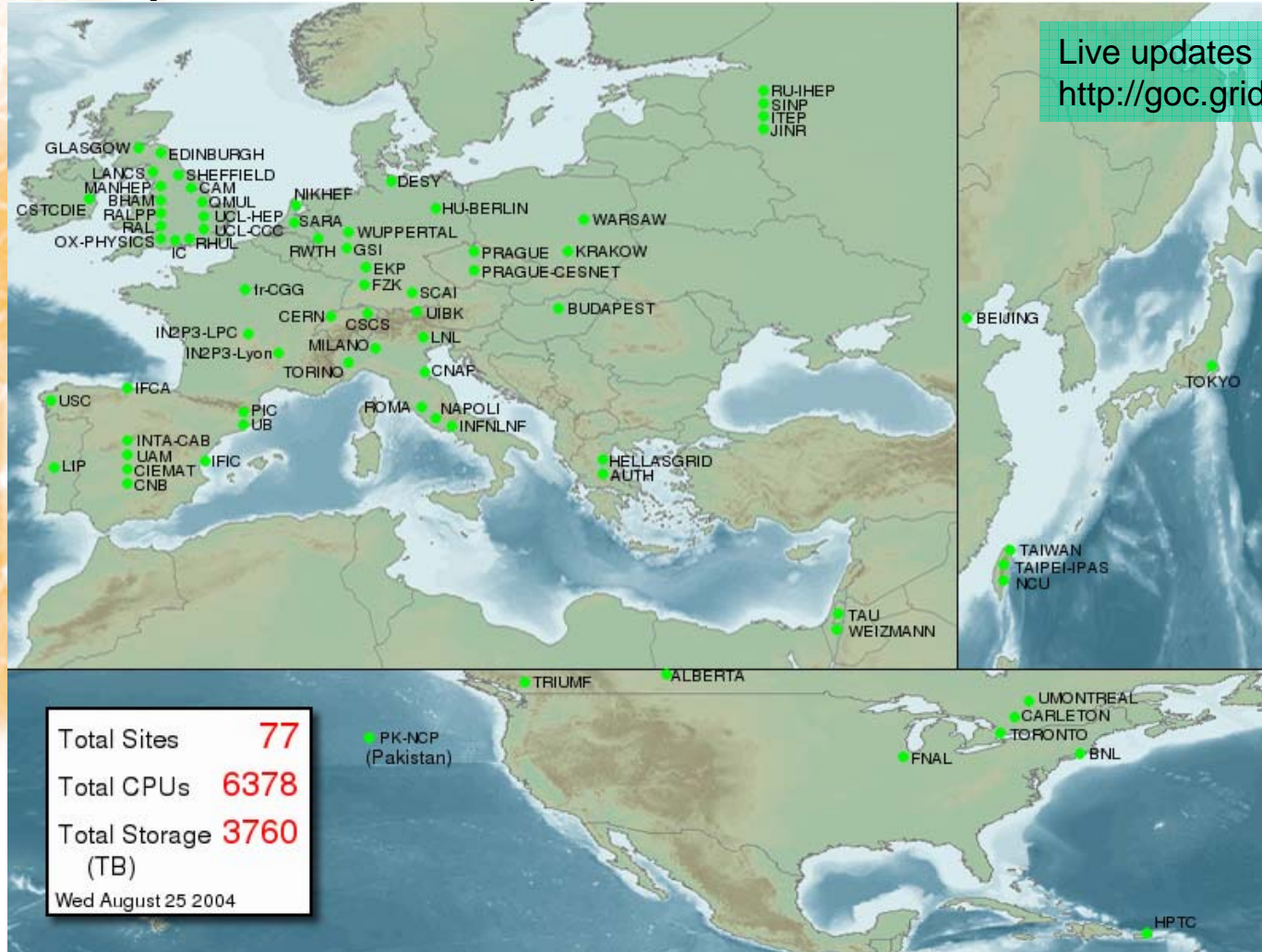  - The research networks

Information Society

**Launched Sept'03 with 12 sites, now more than 100 sites and continues to grow**

Live updates
http://goc.grid-support.ac.uk/lcg2

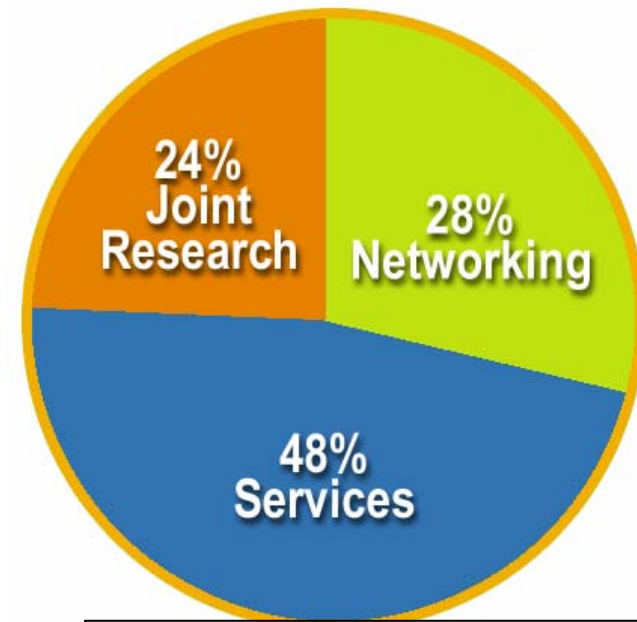| Total Sites | 77 |
| Total CPUs | 6378 |
| Total Storage (TB) | 3760 |

Wed August 25 2004

Information Society

32 Million Euros EU funding over 2 years starting 1st April 2004

- 48 % service activities **(Grid Operations, Support and Management, Network Resource Provision)**

- 24 % middleware re-engineering **(Quality Assurance, Security, Network Services Development)**

- 28 % networking **(Management, Dissemination and Outreach, User Training and Education, Application Identification and Support, Policy and International Cooperation)**



24% Joint Research

28% Networking

48% Services

**Emphasis in EGEE is on operating a production grid and supporting the end-users**

Information Society

- **Frank finds that EGEE is becoming ubiquitous and so he can connect through his University's connection to its NERN**

- **EGEE is focussing on providing a production gird – so it will be available when he needs it.**

- **EGEE is geared towards its users**

- **EGEE - *what is it and why is it needed?***
- Middleware – *current and future*
- **Operations – *providing a stable service***
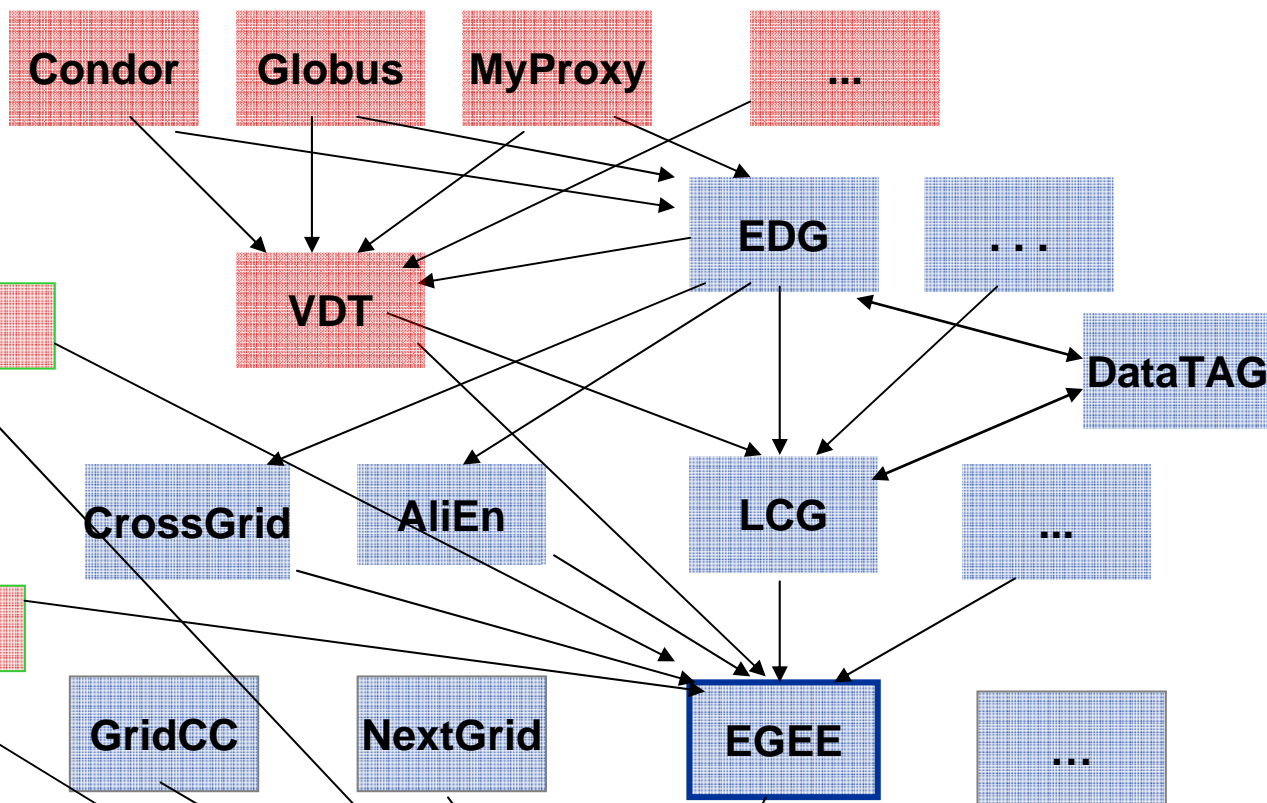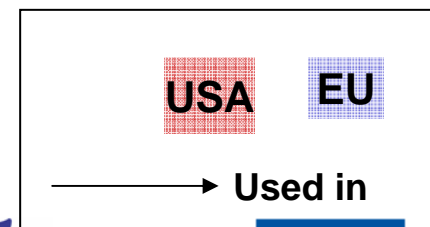- **Networking – *enabling collaboration***
- **Summary**

**Information Society**

- **"gLite" - the new EGEE middleware**
- **Service oriented - components that are :**
  - Loosely coupled (by messages)
  - Accessible across network; modular and self-contained; clean modes of failure
  - So can change implementation without changing interfaces
  - Can be developed in anticipation of new uses
- **… and are based on standards.  Opens EGEE to:**
  - New middleware (plethora of tools now available)
  - Heterogeneous resources (storage, computation…)
  - Interact with other Grids (international, regional and national)

- **Lightweight (existing) services**
  - Easily and quickly deployable
  - Use existing services where possible as basis for re-engineering
- **Interoperability**
  - Allow for multiple implementations
- **Resilience and Fault Tolerance**

- **Co-existence with deployed infrastructure**
  - Reduce requirements on site components
  - Co-existence (and convergence) with LCG-2 and Grid3 are essential for the EGEE Grid service

- **Service oriented approach**
  - Follow WSRF standardization
  - No mature WSRF implementations exist to date so start with plain WS (WS-I)
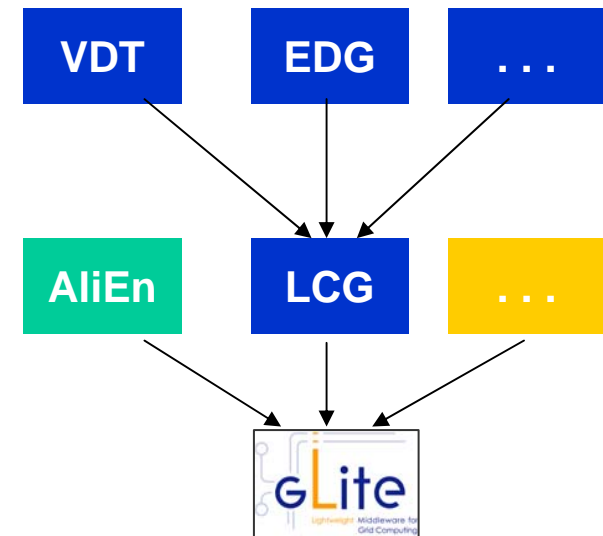  - Provide framework to others so higher-level services can be developed quickly

    Architecture: https://edms.cern.ch/document/476451

Information Society

- **Exploit experience and components from existing projects**
  - AliEn, VDT, EDG, LCG, and others
- **Design team works out architecture and design**
  - Feedback and guidance from EGEE PTF & applications; Operations, LCG GAG & ARDA
- **Components are initially deployed on a prototype infrastructure**
  - Small scale (CERN & Univ. Wisconsin)
  - Get user feedback on service semantics and interfaces
- **After internal integration and testing, components are delivered to grid operations group and deployed on the pre-production service**
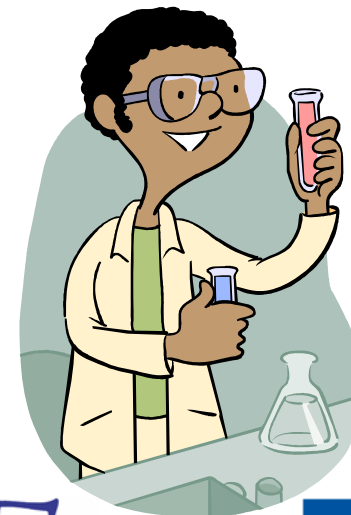
Information Society

- Intended to replace LCG-2

- Starts with existing components from AliEN, EDG, VDT etc.

- Aims to address LCG-2 shortcoming and advanced needs from applications

- Prototyping short development cycles for fast user feedback

- Initial web-services based prototypes being tested with representatives from the application groups

- **EGEE - *what is it and why is it needed?***
- **Middleware – current and future**
- Operations – *providing a stable service*
  - Needs more than middleware
  - Organisational, operational infrastructure
- **Networking – *enabling collaboration***
- **Summary**

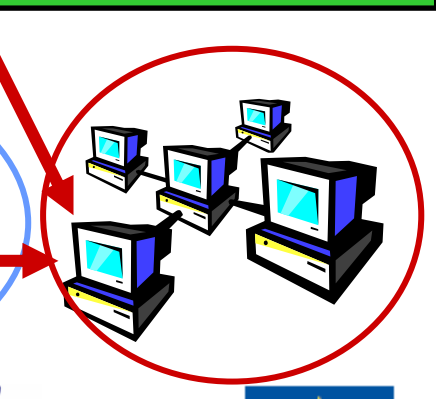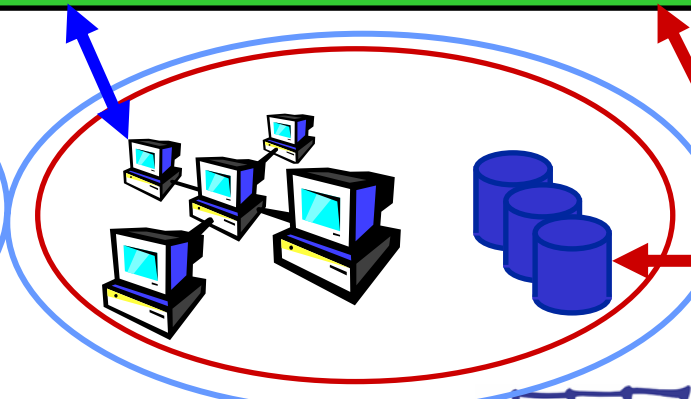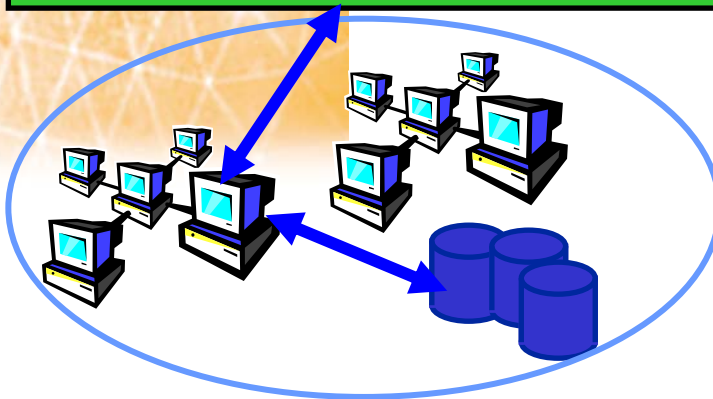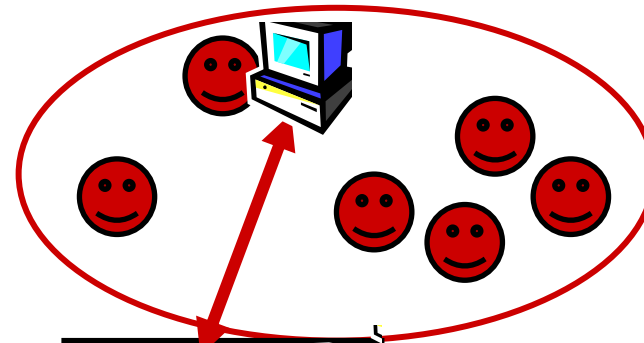- **Frank finds that EGEE is building on existing software – so he can find existing applications and documentation**

- **EGEE is building service oriented middleware, so Frank can bring his understanding of Web Services and can connect biomed services**

Information Society

**User-view of EGEE:**
**a multi-VO Grid**
Enabling Grids for E-sciencE

**Grid services**

User Interface

User Interface

Information Society

**EGEE has a formal procedure for adding selected new user communities (Virtual Organisations):**

- **Negotiation with one of the Regional Operations Centres**

- **Seek balance between the resources contributed by a VO and those that they consume.**

- **Resource allocation will be made at the VO level.**

- **Many resources need to be available to multiple VOs : shared use of resources is fundamental to a Grid**

- Authentication
  - User obtains certificate from CA
  - Connects to UI by ssh
  - Downloads certificate
  - Invokes Proxy server
  - Single logon – to UI - then Secure Socket Layer with proxy identifies user to other nodes

- Authorisation - currently
  - User joins Virtual Organisation
  - VO negotiates access to Grid nodes and resources (CE, SE)
  - Authorisation tested by CE, SE:

  gridmapfile maps user to local account

**CA**

*Personal*

**VO mgr**

**UI**

**VO service**

**VO database**

**SSL**

**(proxy)**

**Gridmapfiles**

**On CE, SE nodes**

Information Society

# Running the Production  Service

**Grid deployment has entered a new phase**

- Basic middleware is working
  - responsible now for a small fraction of the problems
- Outstanding performance/functionality issues
  - RLS, RB /  little modularity & lack of consistent interfaces …
  - some solutions are being developed but many cannot be addressed in current software/architecture - *set priorities for new middleware* (gLite)
- Many operational issues
  - mis-configuration, out of date mware, single points of failure, failover, mgmt interfaces …
  - resources unsuitable for applications needs (e.g. insufficient disk space)
  - slow response by sites to problems (holiday periods, security concerns)
  - new middleware will not help for many of these issues - grid partners must think *Service*

**The grid still does not appear as a single coherent facility**
applications must adapt to the current service to gain maximum profit
but result has been very effective for LHCb - ~3000 concurrent jobs

# EGEE Operations (I): OMC and CIC

- **Operation Management Centre**
  - located at CERN, coordinates operations and management
  - coordinates with other grid projects

- **Core Infrastructure Centres**
  - behave as single organisations
  - operate core services (VO specific and general Grid services)
  - develop new management tools
  - provide support to the Regional Operations Centres

○ Operations Management Centre

● Core Infrastructure Centre

○ Regional Operations Centre

# EGEE Operations (II): ROC

- **Regional Operations Centre  responsibilities and roles:**
  - Testing (certification) of new middleware on a variety of platforms before deployment
  - Deployment of middleware releases + coordination + distribution inside the region
  - integration of 'Local' VO
  - Development of procedures and capabilities to operate the resources
  - First-line user support
  - Bring new resources into the infrastructure and support their operation
  - Coordination of integration of national grid infrastructures Provide resources  for pre-production service

Information Society

- **Frank finds that EGEE has a user support structure which will have a local center for him interact with.**

- **A production grid has high availability – 24/7 and accessible anywhere**

- **EGEE -** *what is it and why is it needed?*

- **Middleware – current and future**

- **Operations – providing a stable service**

- Networking – *enabling collaboration*

  - *Current application communities*

- **Summary**

# EGEE pilot application: Large Hadron Collider



- **Data Challenge:**
  - 10 Petabytes/year of data !!!
  - 20 million CDs each year!

- **Simulation, reconstruction, analysis:**
  - LHC data handling requires computing power equivalent to ~100,000 of today's fastest PC processors!

- **Operational challenges**
  - Reliable and scalable through project lifetime of decades



Mont Blanc
(4810 m)

Downtown Geneva

# EGEE pilot application: BioMedical

- BioMedical
  - Bioinformatics (gene/proteome databases distributions)
  - Medical applications (screening, epidemiology, image databases distribution, etc.)
  - Interactive application (human supervision or simulation)
  - Security/privacy constraints
    - Heterogeneous data formats - Frequent data updates - Complex data sets - Long term archiving
- BioMed applications deployed
- **GATE -** Geant4 Application for Tomographic Emission
  - **GPS@ -** genomic web portal
  - **CDSS -** Clinical Decision Support System

Information Society

- **GATE:** Geant4 Application for Tomographic Emission (LPC)

- **Docking platform for tropical diseases:** grid-enabled docking platform for in sillico drug discovery (LPC)

- **CDSS:** Clinical Decision Support System (UPV)

- **GPS@:** Grid genomic web portal (IBCP)

- **SiMRI 3D:** Magnetic Resonance Image simulator (CREATIS)

- **gPTM 3D:** Interactive radiological image visualization and processing tool (LRI)

- **xmipp_ML_refine:** Macromolecular 3D structure analysis (CNB)

- **xmipp_multiple_CTFs :** Electronmicroscopic images CTF calculation (CNB)

- **GridGRAMM:** Molecular Docking web (CNB)

- **GROCK:** Mass screenings of molecular interaction (CNB

- **Mammogrid:** Mammograms analysis (EU project)

- **SPLATCHE:** Genome evolution modeling (U. Berne/WHO)

- **SPLATCHE**
  - first application being migrated from GILDA to biomed VO
- **Pharmacokinetics in MRI (UPV)**
  - MRI registration for contrast agent diffusion study
- **Some progress on biological sequences analysis (M. Lexa)**
- **...**

# BLAST – comparing DNA or protein sequences

- 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

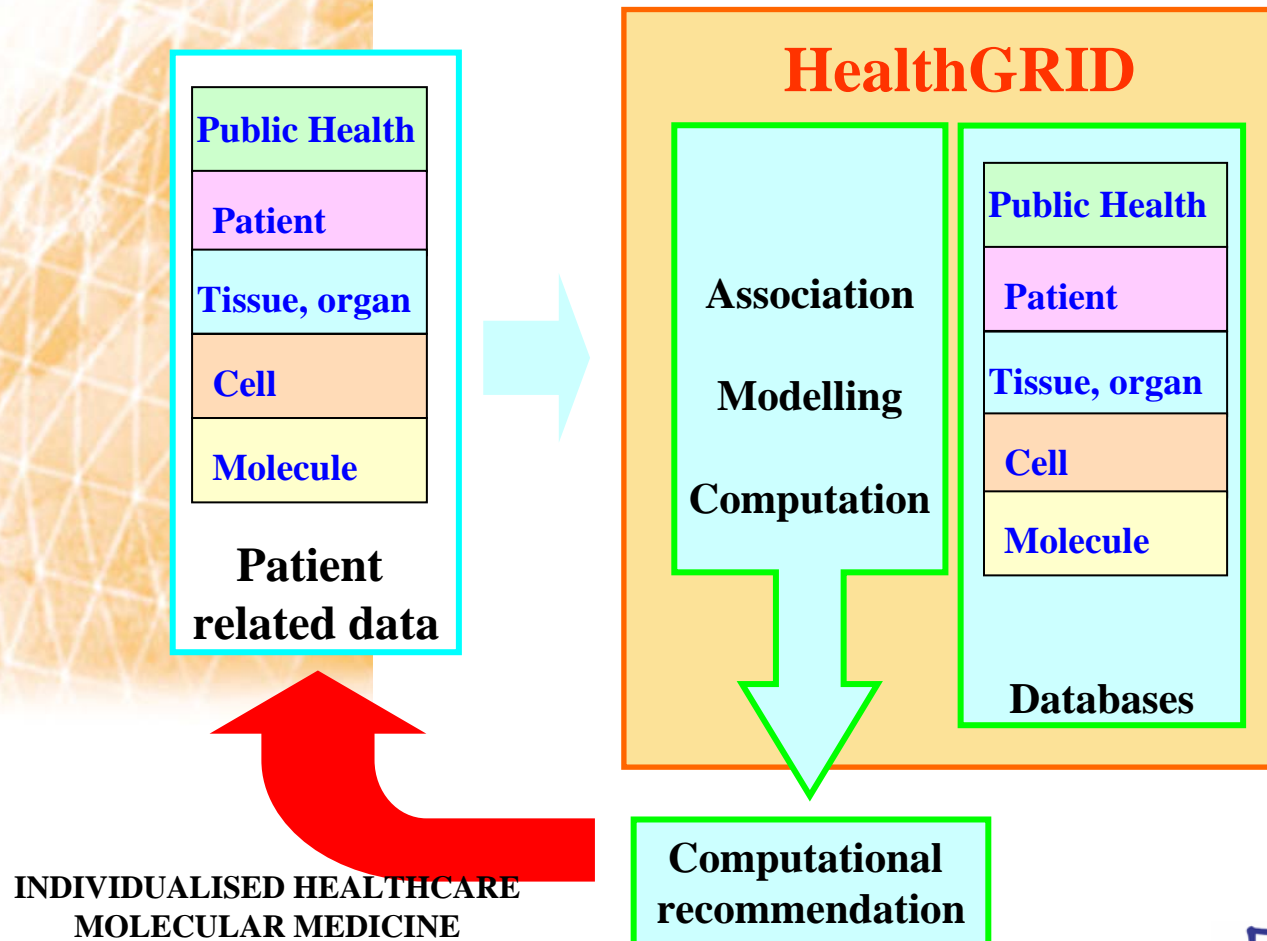www.eu-egee.org

Information Society

- **Frank finds that there are some Biomed applications already available on the Grid.**

- **Some of these may be ones he wants to use.**

- **With a SOA he can compose the services using web service interfaces**

# A look at the future: the HealthGrid vision

In this context "Health" does not involve only clinical practice but covers the whole range of information from molecular level (genetic and proteomic information) over cells and tissues, to the individual and finally the population level (social healthcare).

## HealthGRID

**Public Health**

**Patient**

**Tissue, organ**

**Cell**

**Molecule**

**Patient related data**

Association

Modelling

Computation

**Public Health**

**Patient**

**Tissue, organ**

**Cell**

**Molecule**

**Databases**

**Computational recommendation**

**INDIVIDUALISED HEALTHCARE
MOLECULAR MEDICINE**

# Earth Sciences in EGEE

- **Research**
  - Earth observations by satellite
    - (ESA(IT), KNMI(NL), IPSL(FR), UTV(IT), RIVM(NL),SRON(NL))
  - Climate :
    - DKRZ(GE),IPSL(FR)
  - Solid Earth Physics:
    - IPGP (FR)
  - Hydrology:
    - Neuchâtel University (CH)
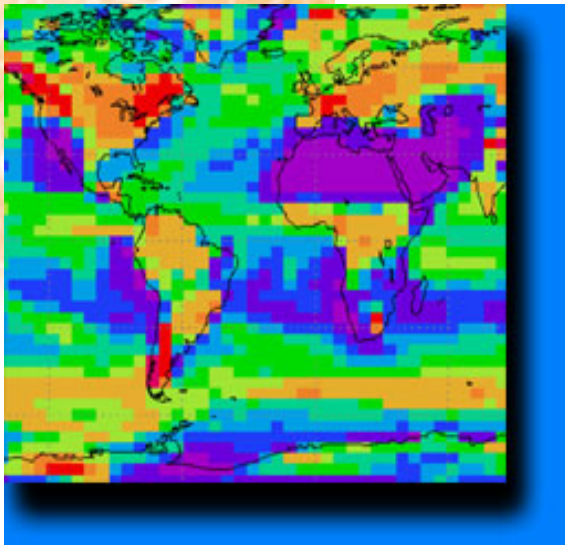- **Industry**
  - CGG : Geophysics Company (FR)

Information Society

**egee**

**Model**: Atmosphere, Ocean, Hydrology, Atmospheric and Marine chemistry….

**Goal:** Comparison of model outputs from different runs and/or institutes

❖Large volume of data (TB) from different model outputs, and experimental data

❖Run made on supercomputer
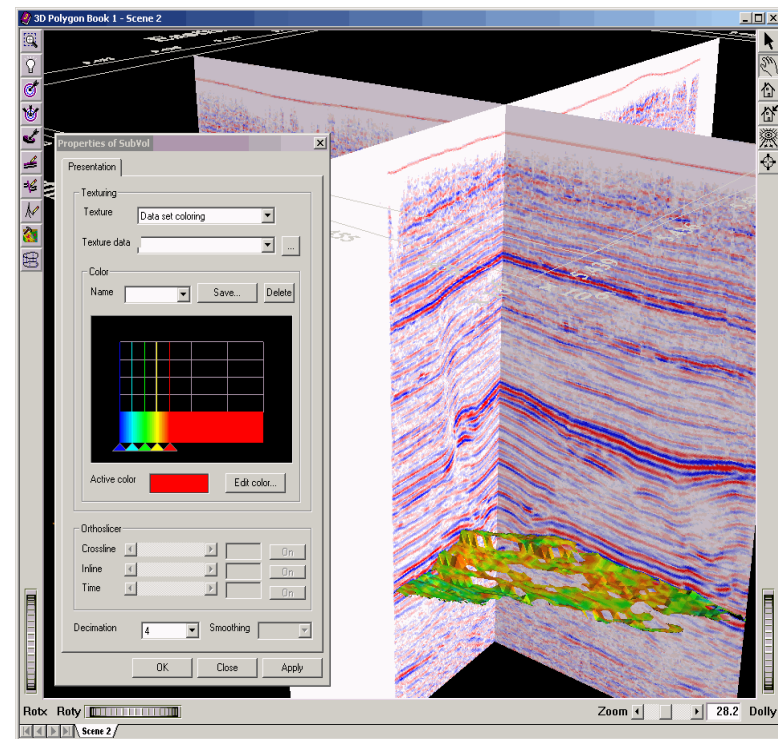=> Link the EGEE infrastruture with supercomputer Grids (DEISA)



EXAMPLE: For the IPCC Assessment reports many experiment are performed with different models (different spatial resolution, different time-step, different "physics" ..) and various sites.
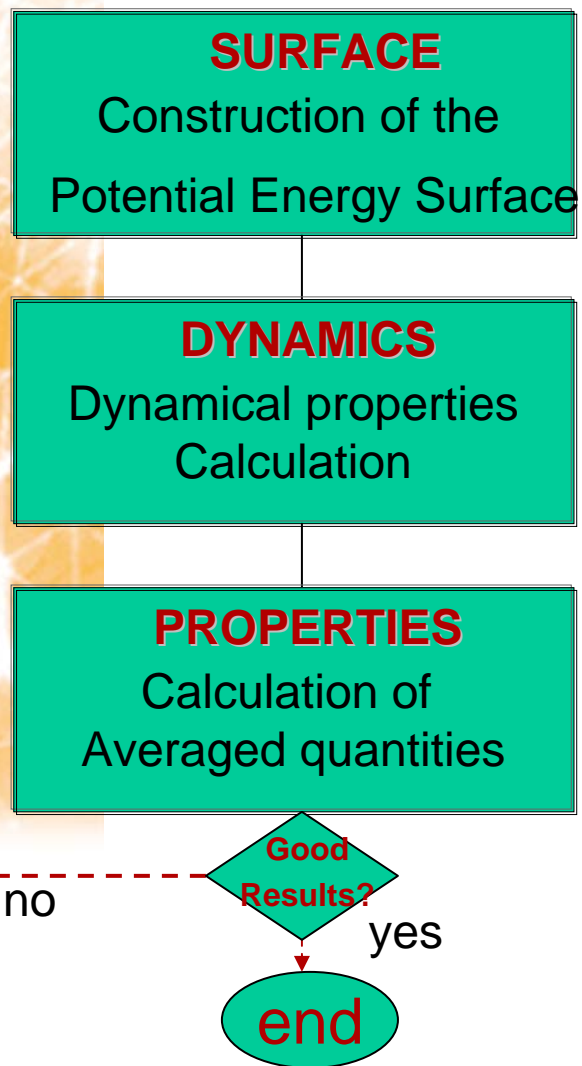
The generated data need to be compared in a comprehensive and "unified" way.

Information Society

**Seismic processing Generic Platform:**

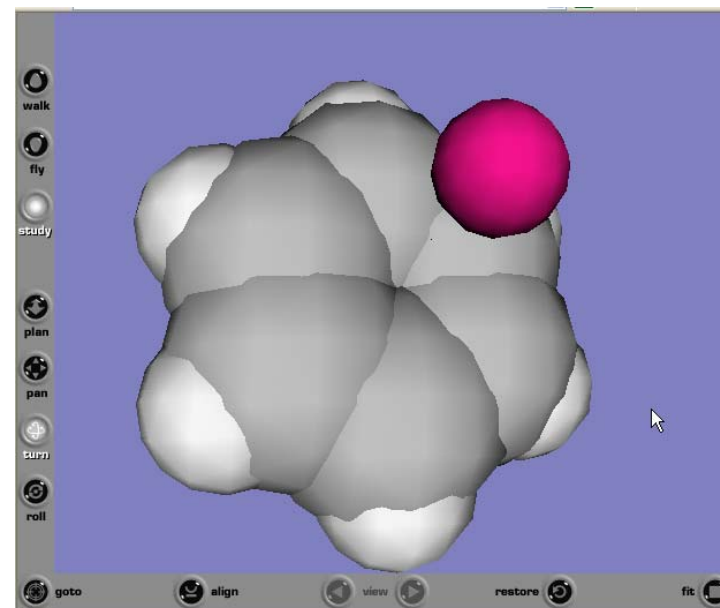- Based on Geocluster, an industrial application – to be a starter of the core member VO.

- Include several standard tools for signal processing, simulation and inversion.

- Opened: any user can write new algorithms in new modules (shared or not)

- Free for academic research

-Controlled by license keys (opportunity to explore license issue at a grid level)

- initial partners F, CH, UK, Russia, Norway

Information Society

**SURFACE**

Construction of the

Potential Energy Surface

**DYNAMICS**

Dynamical properties
Calculation

**PROPERTIES**

Calculation of
Averaged quantities

**Good
Results?**

no

yes

**end**

Ar - Benzene

# Critical Features of the Individual Programs

- **AB INITIO METHODS (molpro, gamess, adc, gaussian, ) resource requests are proportional to** $N^3$ **(N is the number of electrons) and to** $M^D$ **(M is the number of grid points per dimension D) for CPU and disc demand.**

- **EMPIRICAL FORCE FIELDS (Venus, dl_poly, …) resource requests are proportional to P! (P is the number of atoms)**

- DYNAMICS (APH3D, TIMEDEP, …) these programs use as input the output of the previous module most critical dependence is on the total angular momentum **J** value that can increase up to several hundred units and the size of the matrices depend on $2J+1$

- KINETICS PROGRAMS use dynamics results for integrating relevant time dependent applications

# The **MAGIC** telescope

- Largest Imaging Air Cherenkov Telescope **(17 m mirror dish)**

- **Located on Canary Island** La Palma **(@ 2200 m asl)**

- **Lowest** energy threshold **ever obtained with a Cherenkov telescope**

- **Aim: detect** $\gamma$-ray sources **in the unexplored energy range:** 30 *(10)*-> 300 GeV

Information Society

# Applications in EGEE

- **Production service supporting multiple VOs with different requirements**
  - Data
    - Volume
    - Location – distributed?
    - Write Once or Update?
    - Metadata archives?
    - Controlled or open access?
  - Computation
    - High throughput (~ current LCG)
    - High performance, supercomputing
  - No. of sites, scientists,…
- **Establish viable general process to** bring other scientific communities on board

- **EGEE - *what is it and why is it needed?***

- **Middleware – current and future**

- **Operations – providing a stable service**

- Networking – *enabling collaboration*

    - *Current application communities*

    - *Enabling new and effective use of EGEE*
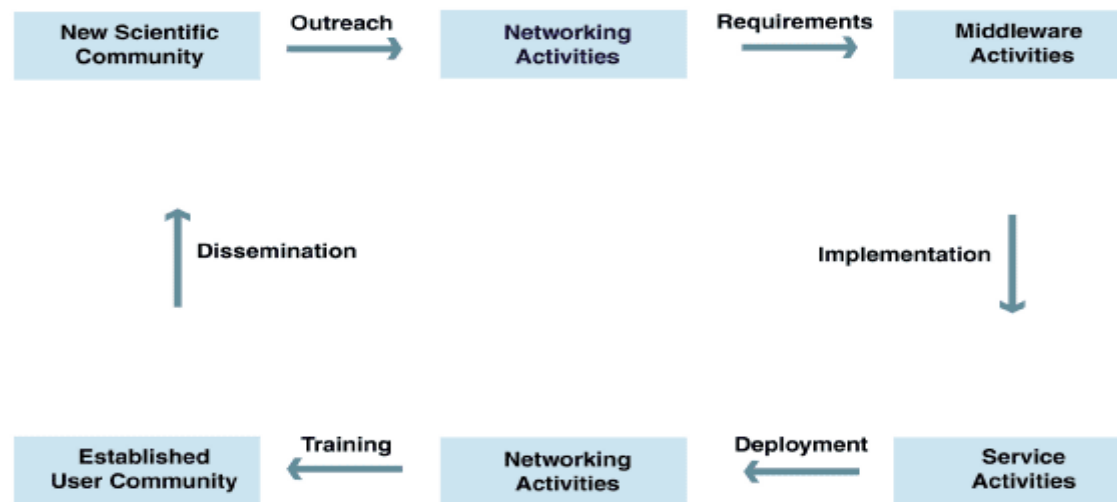
- **Summary**

Information Society

# Who else can benefit from EGEE?

- EGEE Generic Applications Advisory Panel:
  - For new applications

- EU projects: MammoGrid, Diligent, SEE-GRID …

- Expression of interest: Planck/Gaia (astroparticle), SimDat (drug discovery)
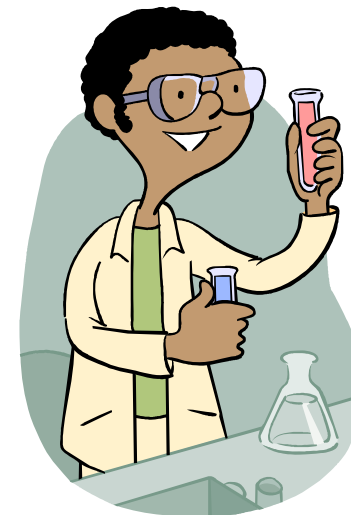
**Information Society**

# Bringing new applications to the grid



1. *Outreach events* inform people about the grid / EGEE
2. Application experts discuss *specific characteristics* with the users
3. *Migrate application* to EGEE infrastructure with the help of EGEE experts
4. *Initial deployment* for testing purposes
5. Production usage - user community contributes computing resources for heavy production demands - "*Canadian dinner party*"

- **Frank finds that there are dedicated activities in EGEE which will help him with bringing his application to the grid**

Information Society

# Intellectual Property

- The existing EGEE grid middleware (LCG-2) is distributed under an Open Source License developed by EU DataGrid
    - Derived from modified BSD - no restriction on usage (academic or commercial) beyond acknowledgement
    - Same approach for new middleware (gLite)

- Application software maintains its own licensing scheme
    - Sites must obtain appropriate licenses before installation

Information Society

- EGEE **is the first attempt to build a worldwide Grid infrastructure for data intensive applications from** many scientific domains

- **A** large-scale production grid service **is already deployed and being used for HEP and BioMed applications with new applications being ported**

- **Resources & user groups will** rapidly expand **during the project**

- **A process is in place for** migrating new applications **to the EGEE infrastructure**

- **A** training programme **has started with events already held**

- **Prototype "*next generation*" middleware is being tested (**gLite**)**

- **Plans for a** follow-on project **are being discussed**