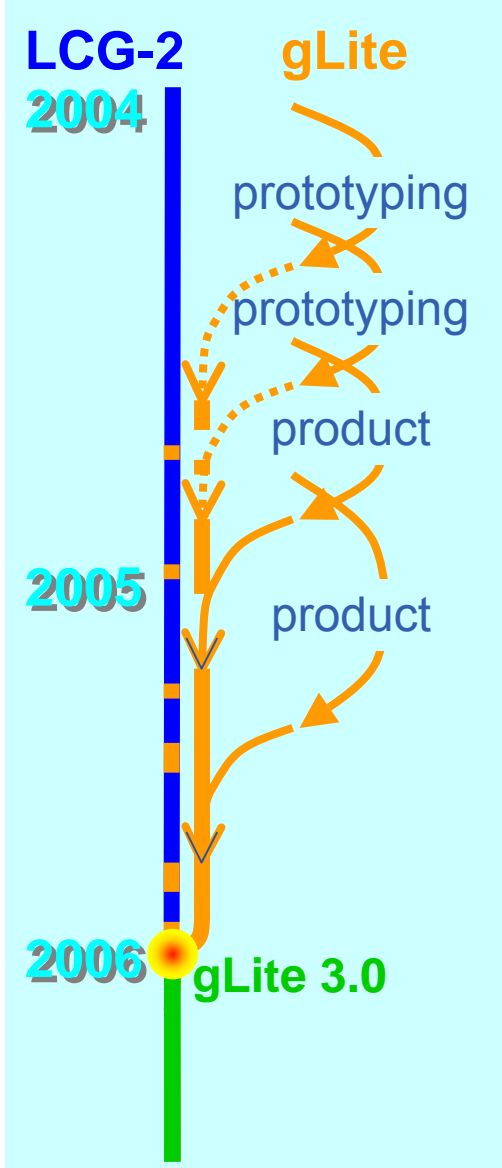# EGEE middleware: gLite

*Claudio Grandi - INFN*

*5th NRENs and Grids Workshop*

*Paris, 11-12 June 2007*
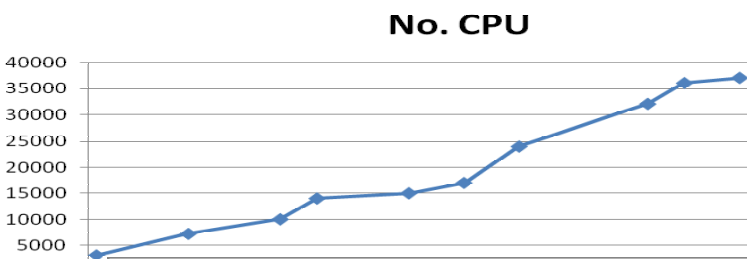
Information Society
and Media

- **Combines components from different providers**
  - Condor and Globus (via VDT)
  - LCG
  - EDG/EGEE
  - Others

- **After prototyping phases in 2004 and 2005 convergence with LCG-2 distribution reached in May 2006**
  - gLite 3.0

- **Focus on providing a deployable MW distribution for EGEE production service**

**LCG-2**   **gLite**

**2004**

prototyping

prototyping

product

**2005**

product

**2006**   **gLite 3.0**

gLite
Lightweight Middleware for Grid Computing

**eGee**

**Enabling Grids for E-sciencE**

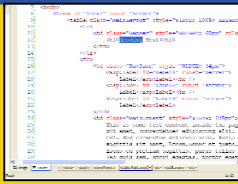**Applications**

**Higher-Level Grid Services**

Workload Management

Replica Management

Visualization

Workflow

Grid Economies

...

**Foundation Grid Middleware**

Security model and infrastructure

Computing (CE) and Storage Elements (SE)

Accounting

Information and Monitoring

- **Applications have access both to Higher-level Grid Services and to Foundation Grid Middleware**

- **Higher-Level Grid Services are supposed to help the users building their computing infrastructure but should not be mandatory**

- **Foundation Grid Middleware will be deployed on the EGEE infrastructure**

  – Must be complete and robust

  – Should allow interoperation with other major grid infrastructures

  – Should not assume the use of Higher-Level Grid Services

**Overview paper http://doc.cern.ch//archive/electronic/egee/tr/egee-tr-2006-001.pdf**

Enabling Grids for E-sciencE

**Directives**

Development

External Software

Software

Error Fixing

**Directives**

Integration

Certification

Pre-Production

GGUS
Global Grid User Support

**Deployment Packages**

**Testbed Deployment**

**Problem**

Production Infrastructure

gLite

**Release**

**Integration Tests**

**Fail**

**Pass**

**Functional Tests**

**Pre-Production Deployment**

**Fail**

**Pass**

**Scalability Tests**

**Fail**

**Installation Guide, Release Notes, etc**

**Pass**

**Enabling Grids for E-sciencE**
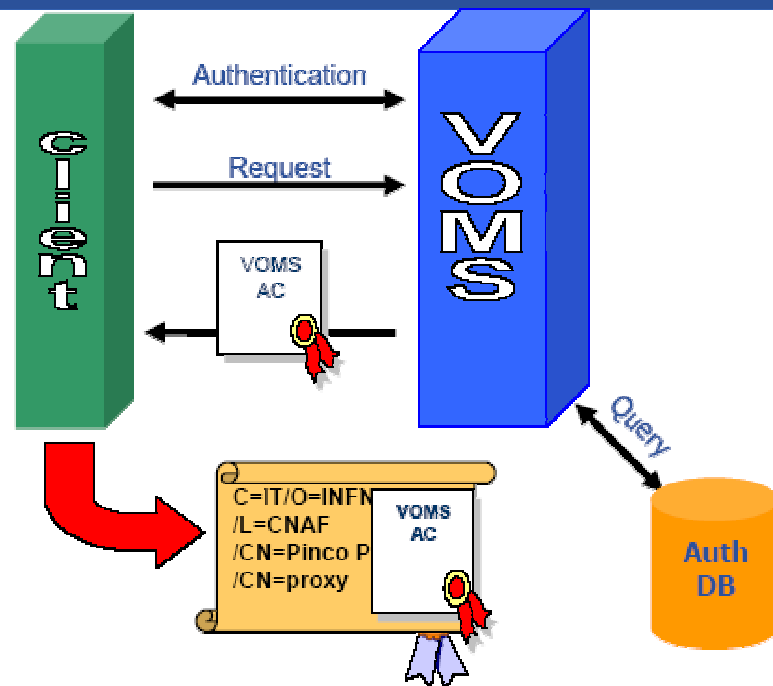
- **Authentication is based on X.509 PKI infrastructure**
  - Certificate Authorities (CA) issue (long lived) certificates identifying individuals (much like a passport)
    - Commonly used in web browsers to authenticate to sites
  - Trust between CAs and sites is established (offline)
  - In order to reduce vulnerability, on the Grid user identification is done by using (short lived) proxies of their certificates
- **Short-Lived Credential Services (SLCS)**
  - issue short lived certificates or proxies to its local users
    - e.g. from Kerberos or from Shibboleth credentials (new in EGEE II)
- **Proxies can**
  - Be delegated to a service such that it can act on the user's behalf
  - Be stored in an external proxy store (MyProxy)
  - Be renewed (in case they are about to expire)
  - Include additional attributes

**Enabling Grids for E-sciencE**

- **VOMS service issues Attribute Certificates that are attached to certificate proxies**
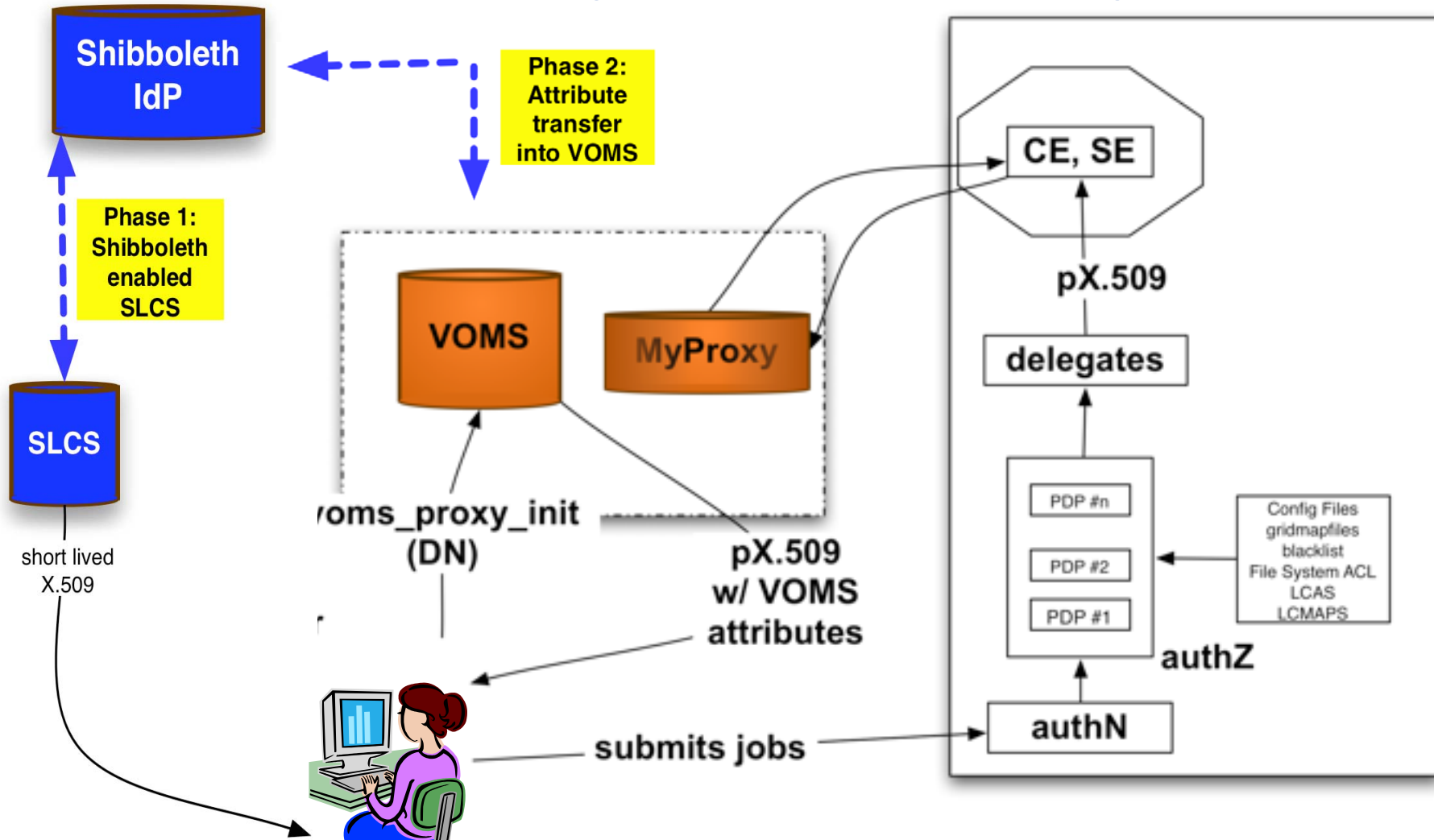  - Provide users with additional capabilities defined by the Virtual Organization
  - Base for the Authorization process
- **Authorization: via mapping to a local user on the resource**
  - glexec changes the local identity (based on suexec from Apache)
  - LCAS/LCMAPS use different plug-ins to determine if and how to map a grid user to a local user
    - mainly used for C-based applications
  - gLite Java Authorization Framework (XACML-compatible)
    - mainly used for Java-based applications
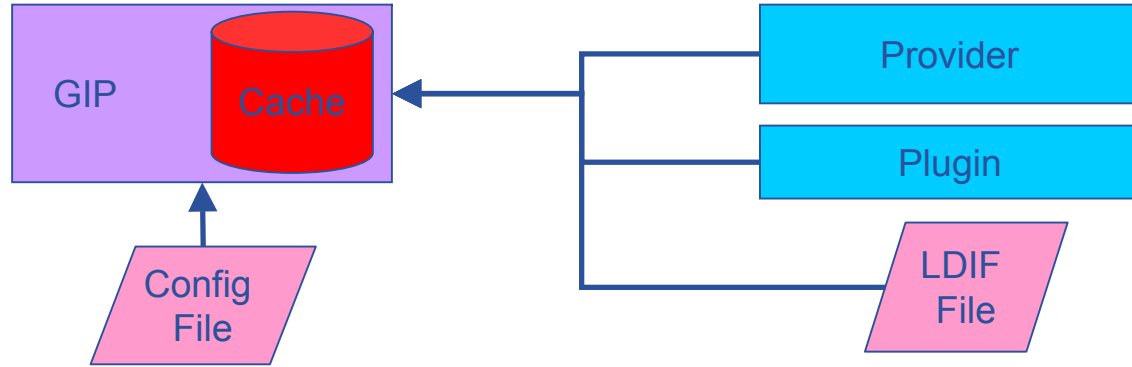  - Compatible with the future *G-PBox* policy management system

**eGee**

Enabling Grids for E-sciencE

## Long lived certificates may be replaced by short lived certificates provided by a Shibboleth identity Provider
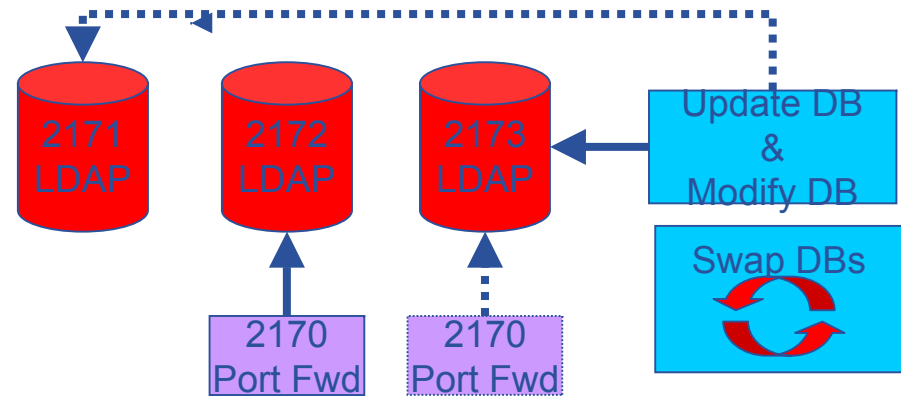
**Enabling Grids for E-sciencE**

- **Generic Information Provider (GIP)**
  - Provides information about a grid service in accordance to the GLUE Schema
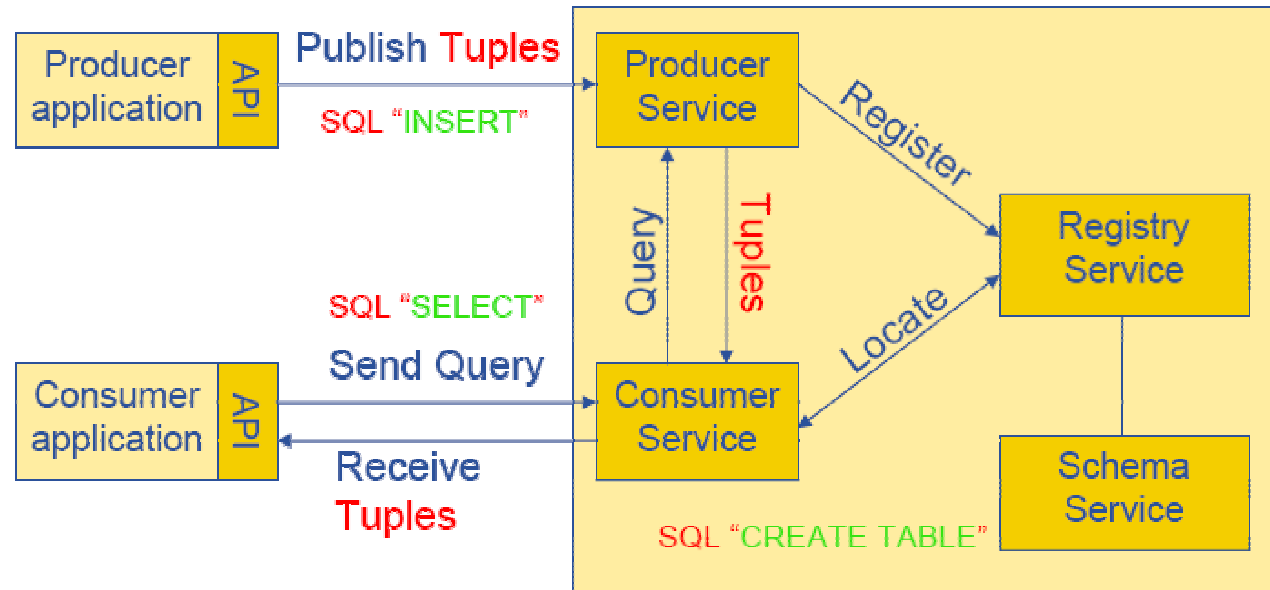- **BDII: Information system**
  - LDAP database that is updated by a process
  - More than one DBs is used separate read and write
  - A port forwarder is used internally to select the correct DB



- **Freedom of choice portal: VOs can white- or black-list resources so that BDII DBs are updated accordingly**
- **Sites failing Site Functional Tests may also be excluded**
- **Up to 2 million queries per day served (over 20 Hz)**

**Enabling Grids for E-sciencE**

- **R-GMA: provides a uniform method to access and publish distributed information and monitoring data**

  – Backbone of EGEE job and infrastructure monitoring

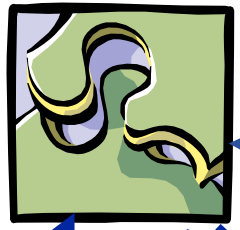  – Working to add authorization



- **Service Discovery: Provides a standard set of methods for locating Grid services**

  – Currently supports R-GMA, BDII and XML files as backends

  – Will add local cache of information
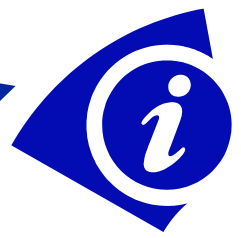
  – Used by some DM and WMS components

User Interface

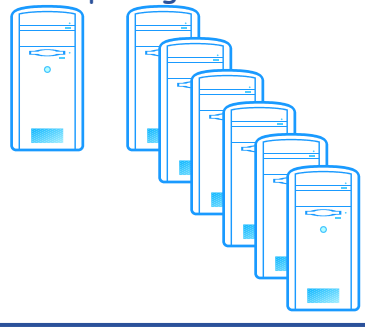Information System

Workload Management
Logging & Bookkeeping

submit

query

retrieve

discover
services

update
credential

publish
state

query

submit

retrieve

publish
state

File and Replica
Catalogs

Site X

Computing Element

Storage Element

Authorization
Service

**Enabling Grids for E-sciencE**
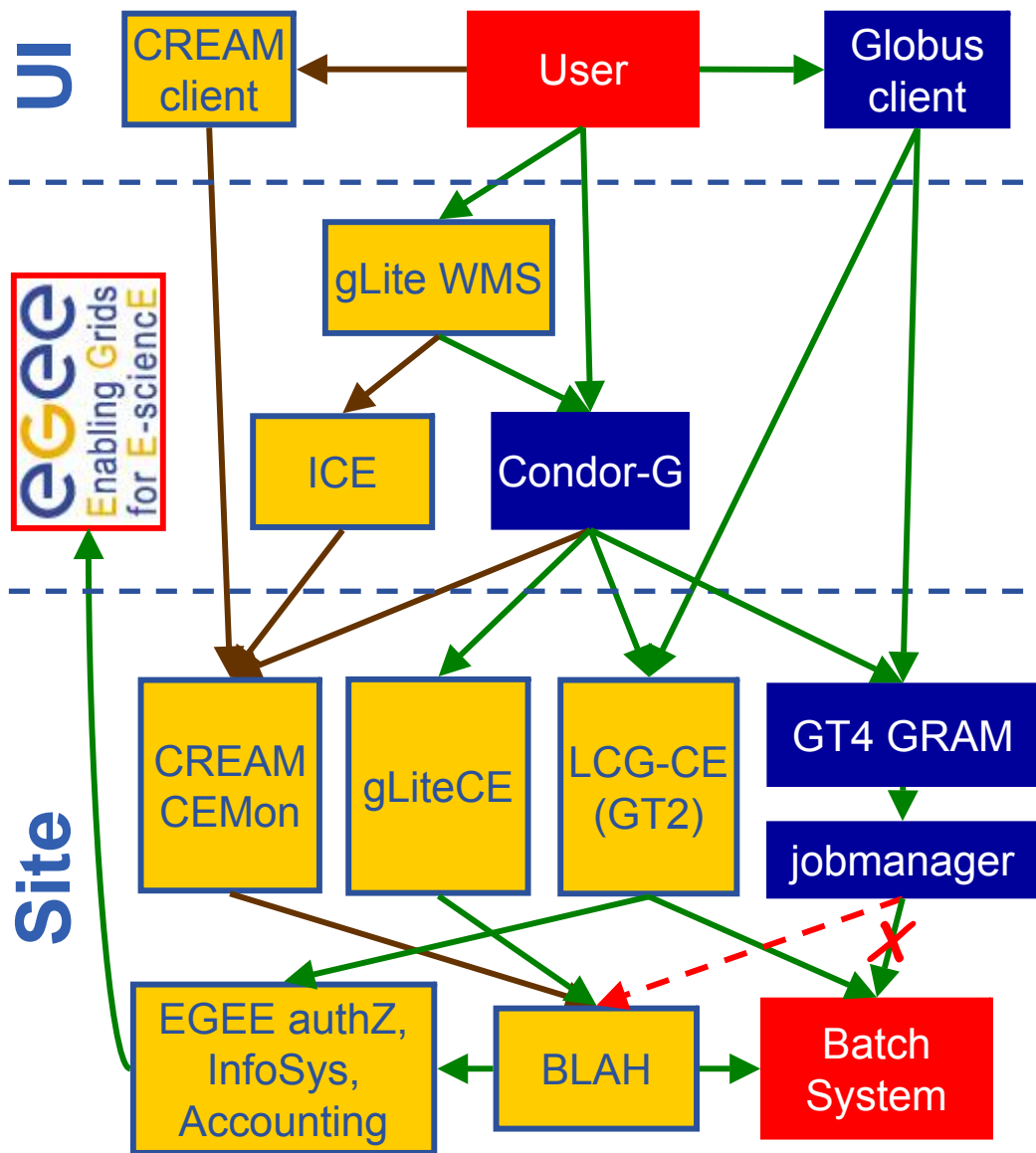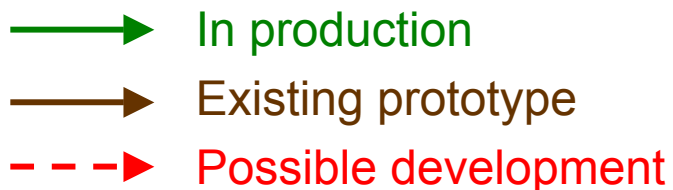
→ **LCG-CE (GT2 GRAM)**
  – Not ported to GT4. To be dismissed

→ **gLite-CE (Condor-C+GSI)**
  – Deployed (GT2 version) but still needs tuning

→ *CREAM (WS-I)*
  – Prototype. OGF-BES (see demo at SC'06)

• **Possible developments:**
  – GT4 → BLAH submissions?

**Choose your preferred path to the Batch System!**

| gLite component | non-gLite component | User / Resource |
|---|---|---|

→ In production

→ Existing prototype

- - -► Possible development

**UI**

CREAM client

User

Globus client

gLite WMS

ICE

Condor-G

**Site**

CREAM CEMon

gLiteCE

LCG-CE (GT2)

GT4 GRAM

jobmanager

EGEE authZ, InfoSys, Accounting

BLAH

Batch System

**Enabling Grids for E-sciencE**
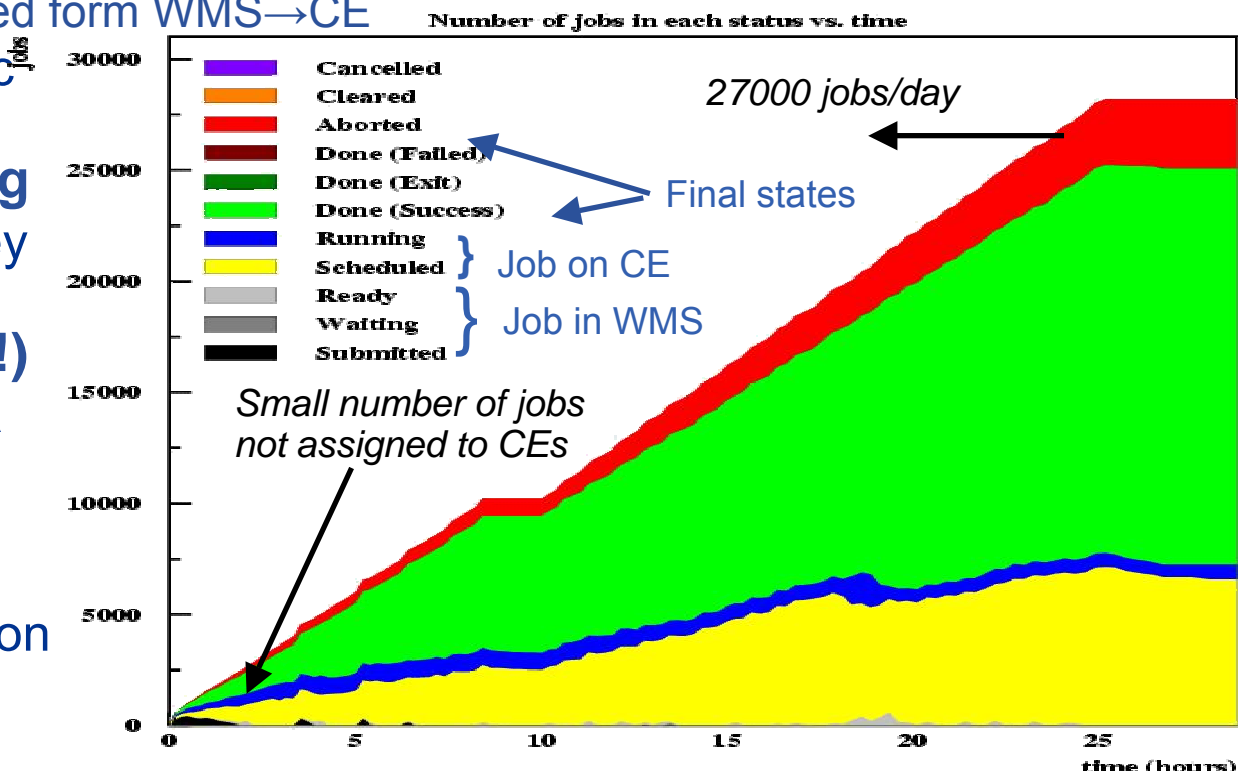
- **Workload Management System**
  - Assigns jobs to resources according to user requirements
    - possibly including data location and user defined ranking of resources
  - Handles I/O data (*input* and *output sandboxes*)
  - Support for compound jobs and workflows (Direct Acyclic Graphs)
    - One shot submission of a group of jobs, shared *input sandbox*
  - Web Service interface: WMProxy
    - UI→WMS decoupled form WMS→CE
  - Support for automatic re-submissions
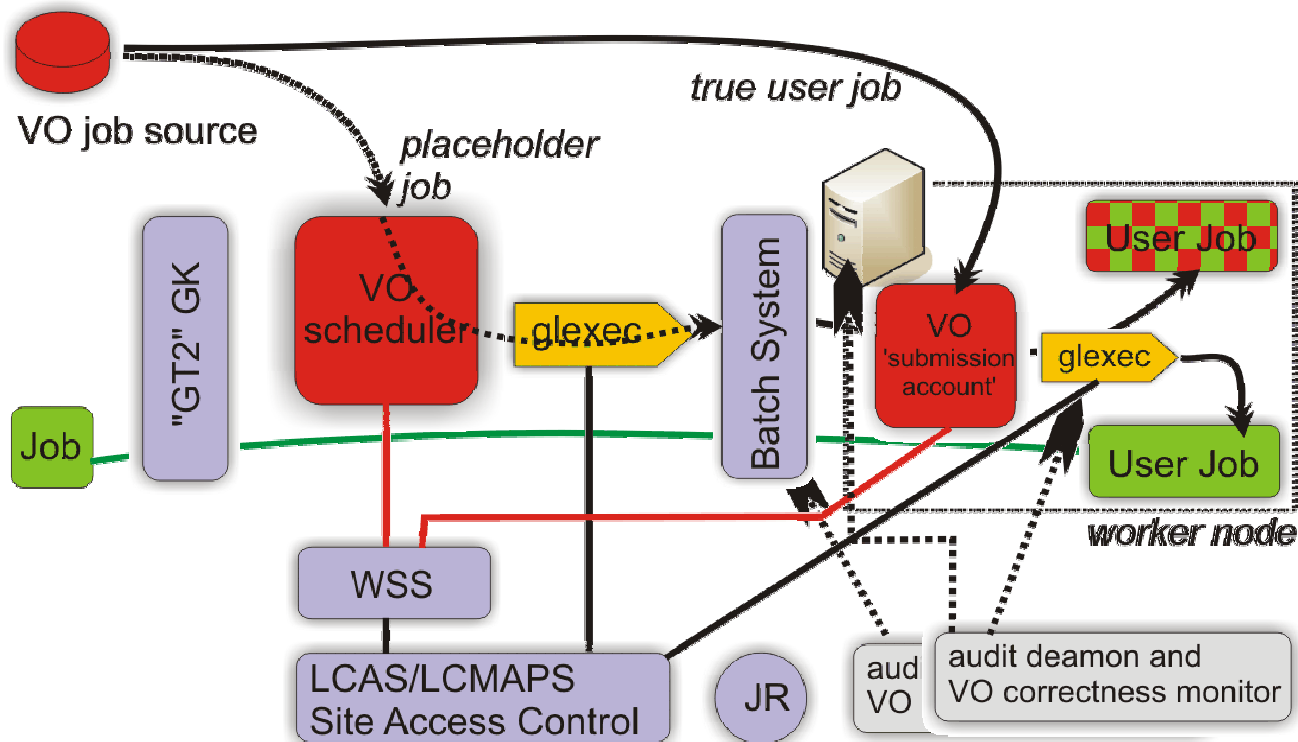- **Logging&Bookkeeping**
  - Tracks jobs while they are running
- **Job Provenance (new!)**
  - Store and retain data on finished jobs
  - Provides data mining capabilities
  - Allows job re-execution

**Number of jobs in each status vs. time**

*27000 jobs/day*

| | |
|---|---|
| Cancelled | |
| Cleared | |
| Aborted | |
| Done (Failed) | |
| Done (Exit) | |
| Done (Success) | Final states |
| Running | } Job on CE |
| Scheduled | |
| Ready | } Job in WMS |
| Waiting | |
| Submitted | |

*Small number of jobs not assigned to CEs*

- **Several VOs submit *pilot jobs* with a single identity for all of the VO**
  - The pilot job gets the user job when it arrives on the WN and executes it
    - Just-in-time scheduling. VO policies implemented at the central queue



- **Use the same mechanism for changing the identity on the Computing Element also on the Worker Nodes (glexec)**
  - The site may know the identity of the real user

**Enabling Grids for E-sciencE**

- **Resource usage by VO, group or single user**
  - *S*ensors running on resources to determine usage
  - It would be possible to enable *Pricing policies* associate a cost to resource usage
    - market-based resource brokering
  - privacy: access to accounting data granted only to authorized people (user, provider, VO manager)



- **Information collected at the Grid Operations Centre (GOC)**
- **Basic functionality in APEL, full functionality in DGAS**

**Enabling Grids for E-sciencE**

- **LFC maps LFNs to SURLs**
  - *Logical File Name* (LFN): user file name
    - in VO namespace, aliases supported
  - *GIbally Unique IDentifier* (GUID)
    - unique string assigned by the system to the file
  - *Site URL* (SURL): identifies a replica
  - A Storage Element and the logical name of the file inside it
- **GSI security: ACLs (based on VOMS)**
  - To each VOMS group/role corresponds a virtual group identifier
  - Support for secondary groups
- **Web Service query interface: Data Location Interface (DLI)**
- **Hierarchical Namespace**
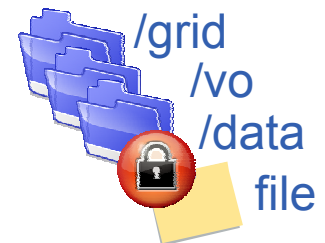- **Supports sessions and bulk operations**

**LFC**

LFN 1 → GUID → SURL 1
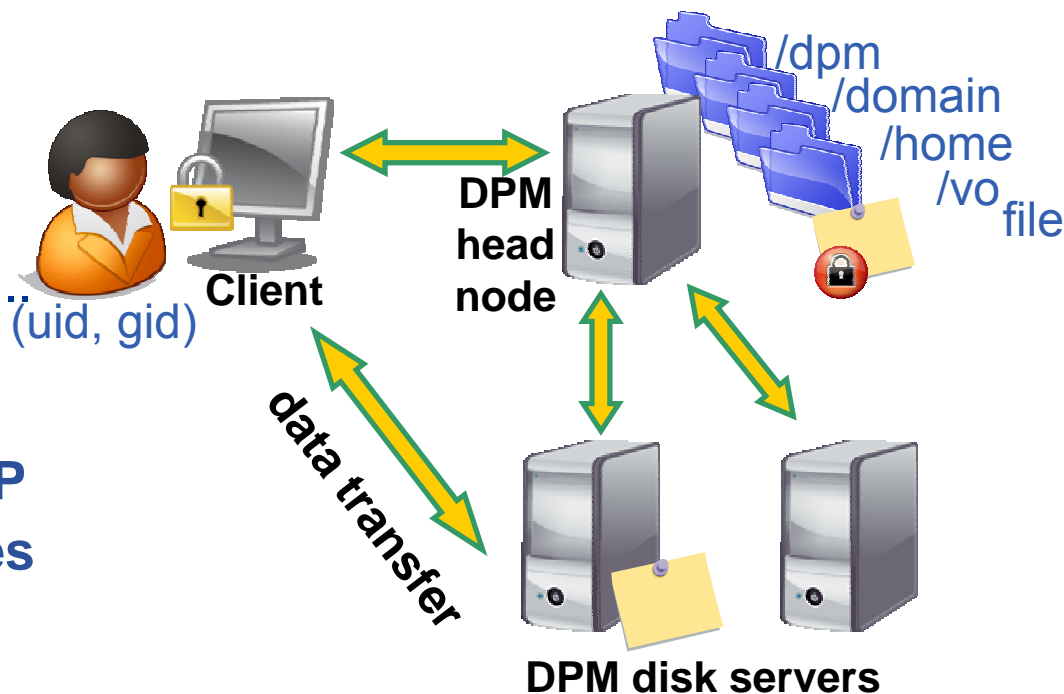LFN 2 → GUID → SURL 2
GUID → ACL

lfc-ls –l /grid/vo/

lfc-getacl /grid/vo/data

**LFC**
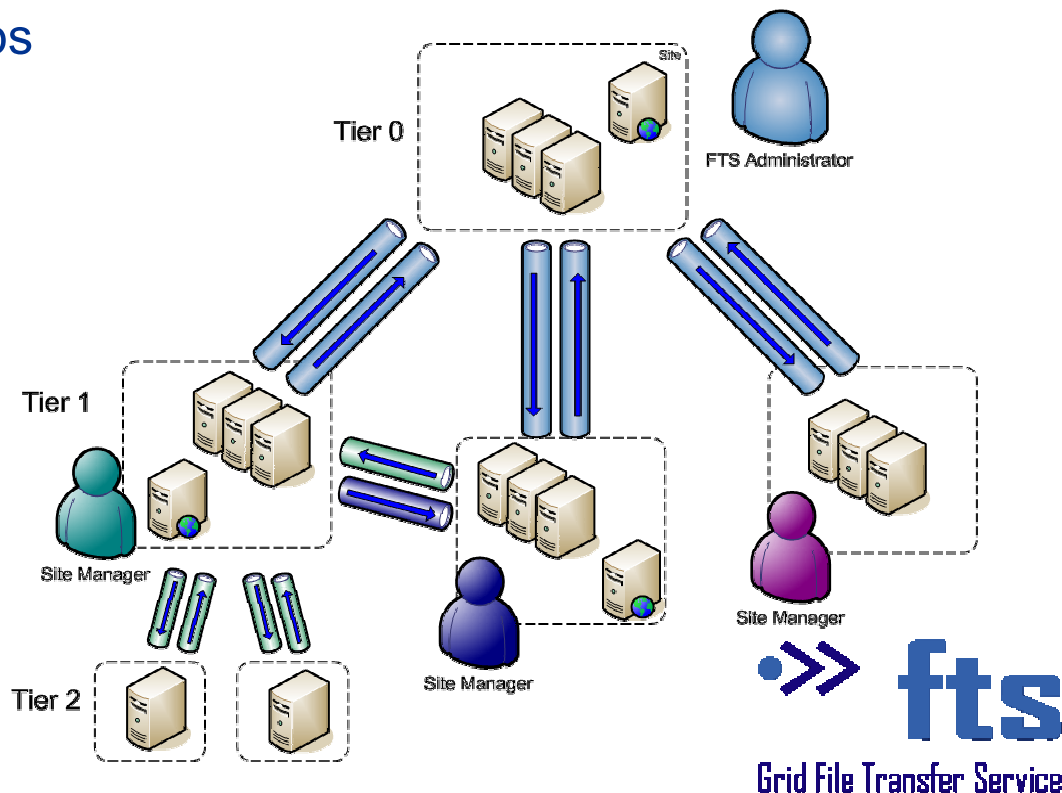**DLI**

/grid
/vo
/data
file

**Enabling Grids for E-sciencE**

- **Storage Resource Manager (SRM): translates SURLs to TURLs**
  - *Transfer URL* (TURL): allows direct access to the file
  - Interface that hides the storage system implementation
  - Handles the authorization based on VOMS credentials
- **Disk-based: DPM, dCache,+; tape-based: Castor, dCache**
- **File I/O: *posix-like* access from local nodes or the grid**
  - GFAL (Grid File Access Layer)
- **Disk Pool Manager (DPM)**
  - Manages storage on disk servers
- **Uses LFS as local catalog**
  - Same features for ACLs, etc...
- **Direct data transfer from/to disk server (no bottleneck)**
- **External transfers via gridFTP**
- **Target: small to medium sites**
  - One or more disk servers

**Client**

**(uid, gid)**

**data transfer**

**DPM head node**

/dpm
/domain
/home
/vo
file

**DPM disk servers**

**Enabling Grids for E-sciencE**

- **FTS: Reliable, scalable and customizable file transfer**
  - Multi-VO service, used to balance usage of site resources according to the SLAs agreed between a site and the VOs it supports
  - WS interface, support for different user and administrative roles (VOMS)
  - Manages transfers through <u>channels</u>
    - mono-directional network pipes between two sites
  - File transfers handled as jobs
    - Prioritization
    - Retries in case of failures
  - Automatic discovery of services
- **Designed to scale up to the transfer needs of very data intensive applications**
  - Demonstrated about 1 GB/s sustained
  - Over 9 petabytes transferred in the last 6 months (> 10 million files)



Grid File Transfer Service
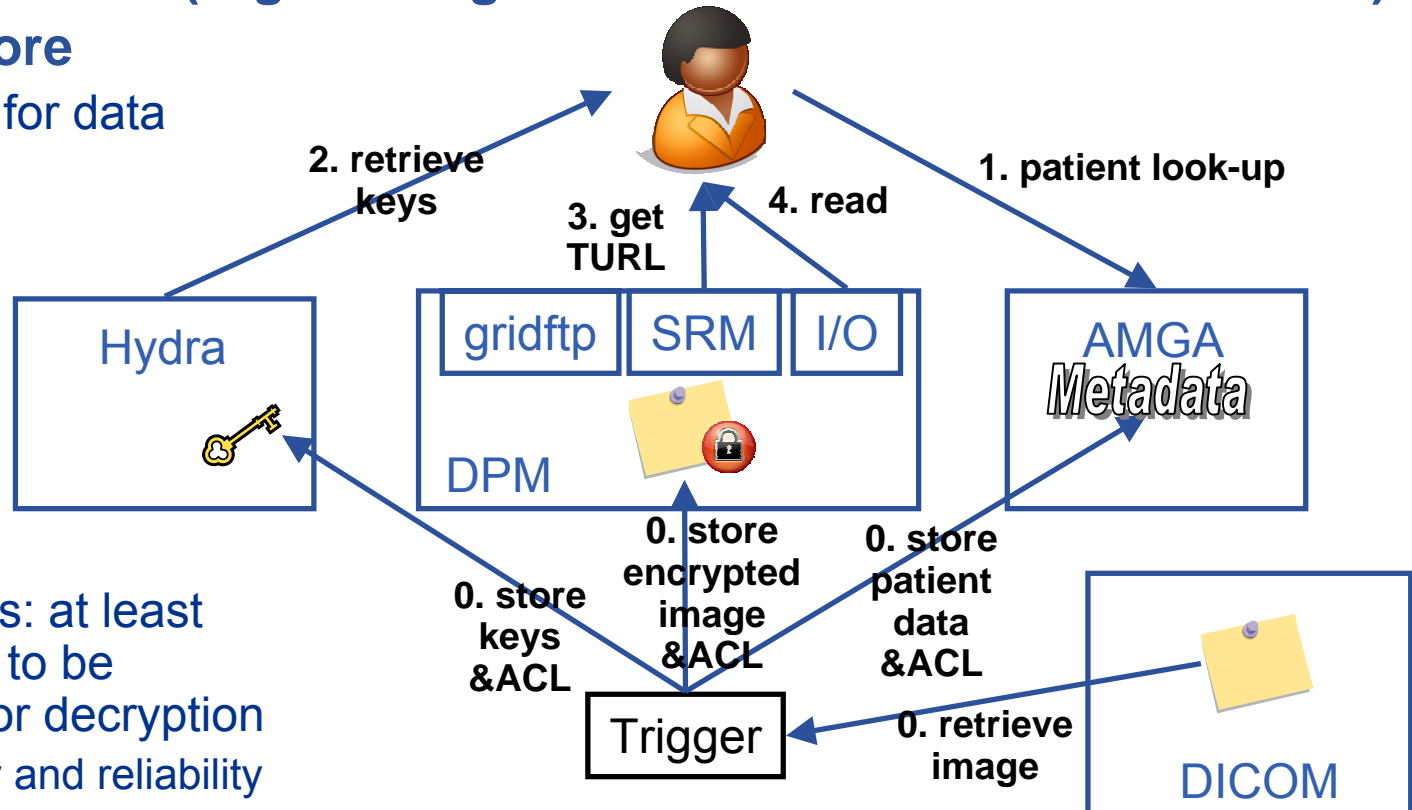
**Enabling Grids for E-sciencE**

- **AMGA is a general purpose metadata catalog**
  - Keeps information about data stored in files
  - Used by several application domains
  - SOAP interface
  - VOMS authorization
  - Shell-like client
  - Graphical Browser (Python)
- **Performance comparable to direct DB access**
  - C++, TCP streaming protocol, very fast SSL sessions
- **LHCb (HEP VO use case)**
  - 100 Million entrie
  - 150GB data
  - $10^5$ entries/day insert rate
  - 10 entries/sec read-rate

**Enabling Grids for E-sciencE**

- **Intended for VO's with very strong security requirements**
  - e.g. medical community
    - anonymity (patient data is separate)
    - fine grained access control (only selected individuals)
    - privacy (even storage administrator cannot read)
- **Interface to DICOM (Digital Image and COmmunication in Medicine)**
- **Hydra keystore**
  - store keys for data encryption

  - N instances: at least M<N need to be available for decryption
    - security and reliability

- **gLite process driven by application and operational requirements**
  - New components added based on their requests and overall importance

- **RESPECT – Program to collect useful tools that work with gLite**
  - See EGEE application portal: http://egeena4.lal.in2p3.fr/index.php
  - Under construction