# Red española de e-Ciencia

## Meta Scheduling and Advanced Application Support on the Spanish NGI

### Enol Fernández del Castillo (IFCA-CSIC)

**OGF 25/EGEE User Forum**

**Catania, March 2nd 2009**

# Outline

1. **Spanish NGI**

2. **Metaschedulers**
   - **Grid-Way**
   - **CrossBroker**

3. **Advanced Application Support**
   - **Interactivity**
   - **MPI**

4. **Summary**

# Spanish NGI: NGI-ES

- National Level entity which operates a general purpose e-science infrastructure
- Objectives:
  - Establish a collaboration framework between all participating institutions to foster a coordinated development of a Grid infrastructure in Spain
  - Propose a sustainable design of the Grid infrastructure that covers the ecosystem of different Grid projects, computing centers, grid infrastructures, etc...
  - Run central services to keep up the infrastructure

Red Española de
e - Ciencia

# Components of NGI-ES

- **Core: Spanish institutions participating in Grid research and development projects:**
  - **EGEE, EUFORIA, DORII, EELA, i2g**, ... with a common middleware based on gLite
  - The new infrastructure **GRID – CSIC  http://www.grid.csic.es**
  - **Universities, computing centers** with Globus Toolkit 4 middleware
  - **RedIris** (Spanish NReN) support EUGRIDPMA certificates.
- Relation with the **Spanish Supercomputing Network (RES)**
  - This network comprises several Spanish research centers that operate a common infrastructure of supercomputing.
  - Analyze possibility of mixed workflows between both infrastructures
- Relation with the **Portuguese NGI in the framework of Ibergrid**

# NGI-ES Interoperability

- Interoperability is needed for the sustainability of the NGI
  - Key issue for the creation of the EGI infrastructure
  - Allows users to select the resources that better fit their necessities
  - Potential access to a significantly larger set of resources
- Reduced management overheads if only a single Grid middleware system needs deployment on each site

# NGI-ES Architecture

## VO Oriented

The architecture of NGI-ES is oriented towards the support of Virtual Organizations

**Key Issues**

- ✓ Advanced VO Services
    - ✓ User Support
    - ✓ Monitoring  & Accounting

- ✓ Application porting and support
- ✓ Middleware driven by applications requirements

**Virtual Organizations**
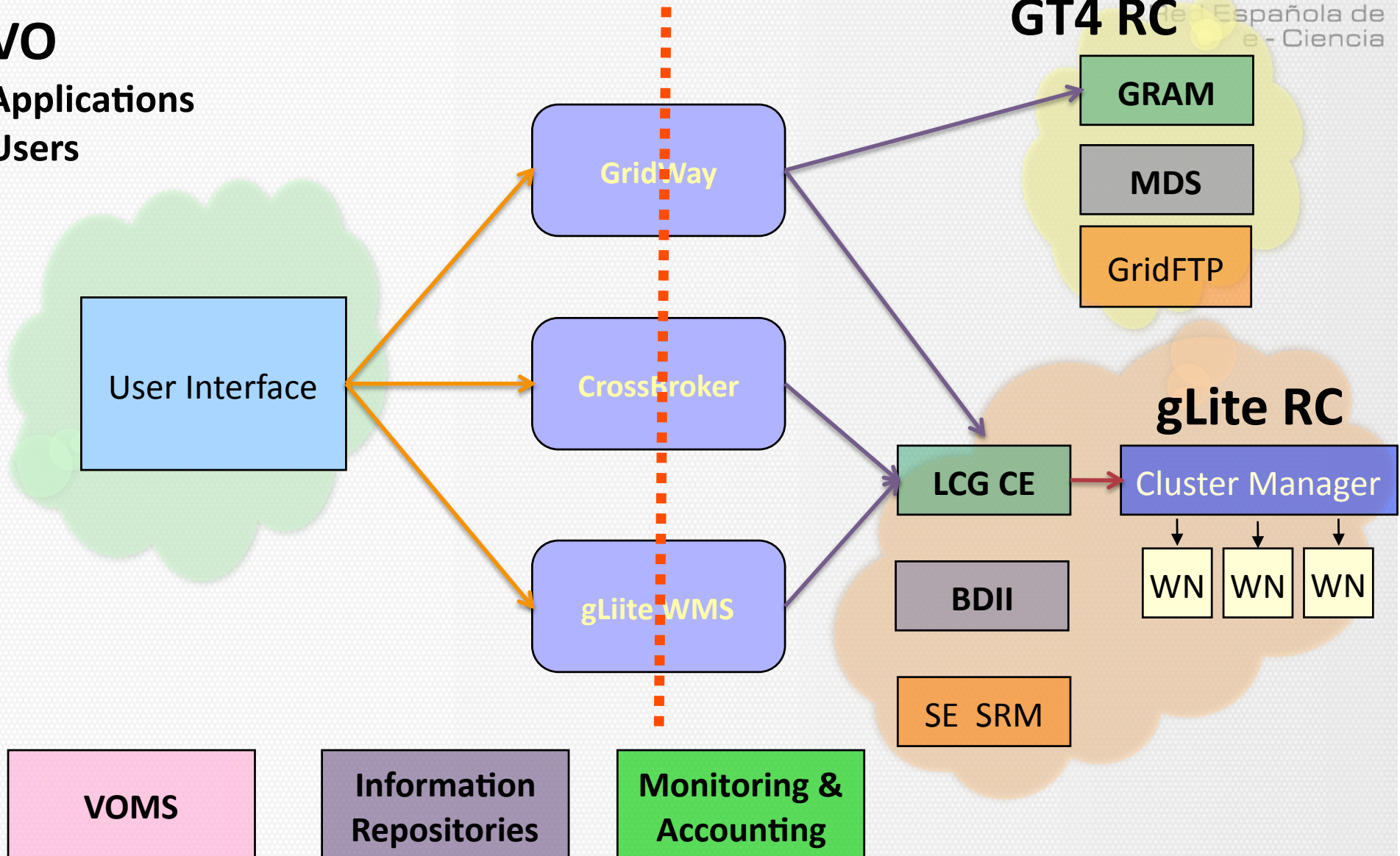
**Central Services**

**RESOURCE CENTERS**

Red Española de
e - Ciencia

gLite

the globus® toolkit

# NGI-ES Architecture

**VO**
**Applications**
**Users**

**GT4 RC**

GRAM

MDS

GridFTP

User Interface

GridWay

CrossBroker

gLiite WMS

**gLite RC**

LCG CE

Cluster Manager

WN  WN  WN

BDII

SE  SRM

**VOMS**
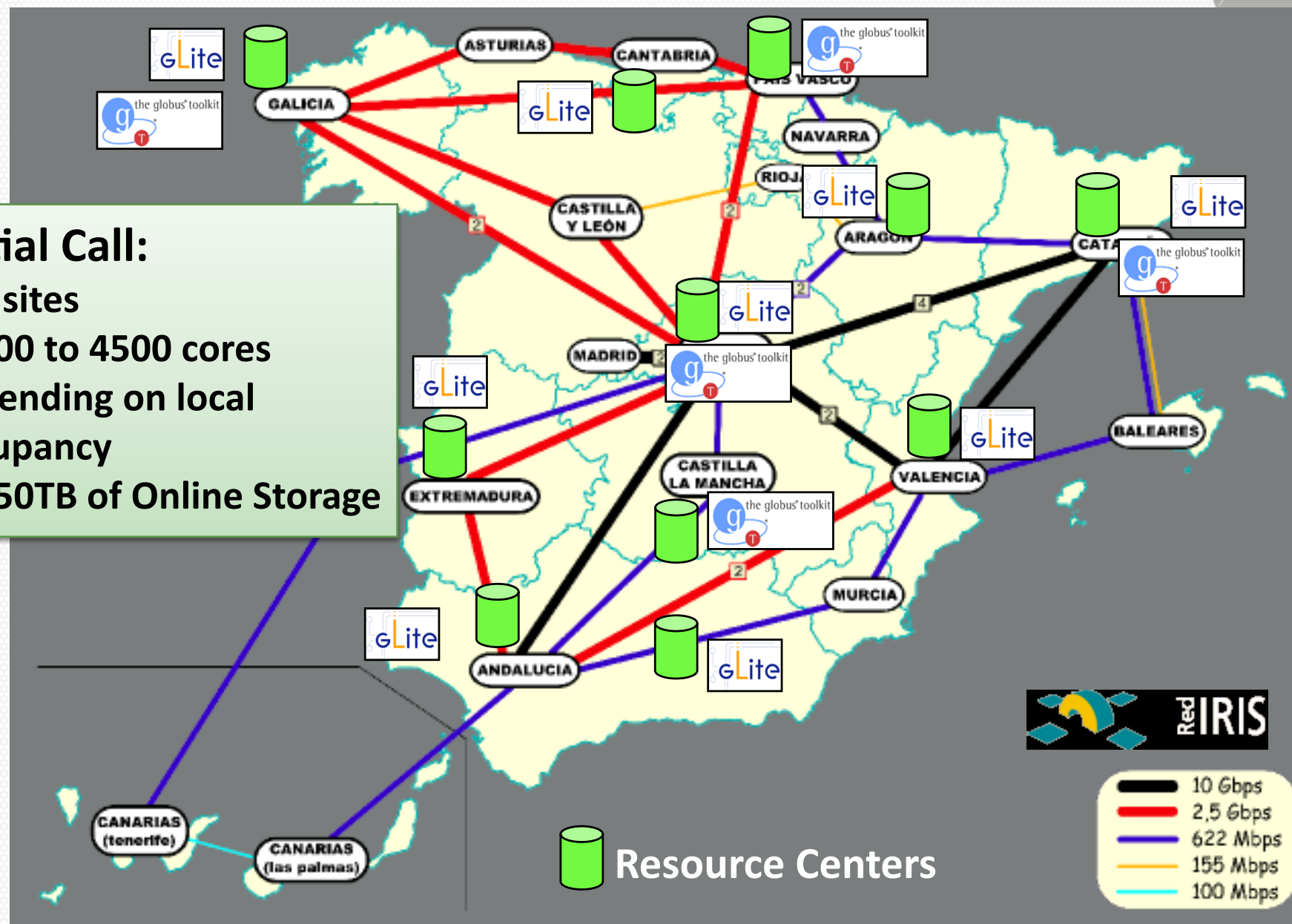
**Information Repositories**

**Monitoring & Accounting**

# Resource Centers Map

**Initial Call:**
- **18 sites**
- **1300 to 4500 cores depending on local occupancy**
- **~350TB of Online Storage**

ASTURIAS CANTABRIA
PAIS VASCO
GALICIA
NAVARRA
RIOJA
CASTILLA Y LEÓN
ARAGON
CATA
BALEARES
MADRID
EXTREMADURA
CASTILLA LA MANCHA
VALENCIA
ANDALUCIA
MURCIA
CANARIAS (tenerife)
CANARIAS (las palmas)

gLite
the globus toolkit

**Resource Centers**

Red IRIS

| | |
|---|---|
| ■ | 10 Gbps |
| ■ | 2,5 Gbps |
| ■ | 622 Mbps |
| ■ | 155 Mbps |
| ■ | 100 Mbps |

# Deployed Services

- **Monitorization and accounting services at Cesga**
  - Accounting Portal: http://www.ngi.cesga.es/gridsite/accounting
  - Monitorization Portal: http://rnagios.ngi.cesga.es/nagios

- **Global Information Repositories:**
  - OpenLDAP + GlueSchema server at IFCA, integrated with NGI-PT
  - Global MDS server at RedIRIS for GT4 resources

- **Metaschedulers:**
  - *GridWay* for the NGI at RedIRIS
  - *Crossbroker* for gLite resources at IFCA
  - *gLite WMS* for bulk submission of serial jobs at IFIC

- **VOMS server** at IFCA

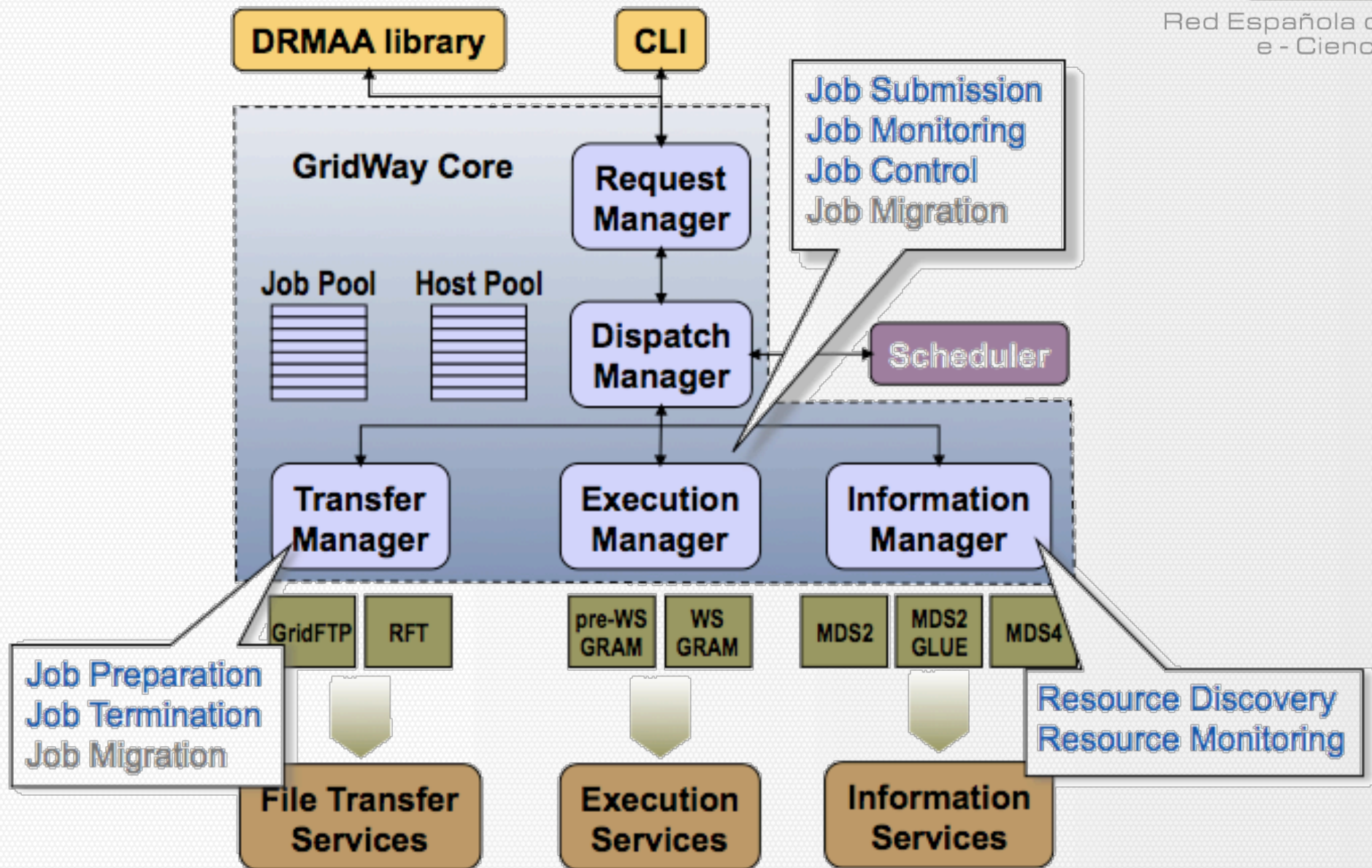- **HelpDesk** at RedIRIS

# NGI-ES Added Value

- Interoperability:
  - Resources from EU projects based on gLite (EGEE, I2G, EELA, WLCG, DORII)
  - Globus Toolkit 4 resources
- Metaschedulers developed by Spanish groups
  - GridWay
  - CrossBroker
- Advanced Application Support
  - Interactivity
  - MPI

# **Metaschedulers: GridWay**

- GridWay allows the efficient use of computing resources of a Grid
  - Included in the Globus distribution
  - Support for both gLite and GT4 resource
- Used by several grid projects and initiatives worldwide

- Developed by UCM
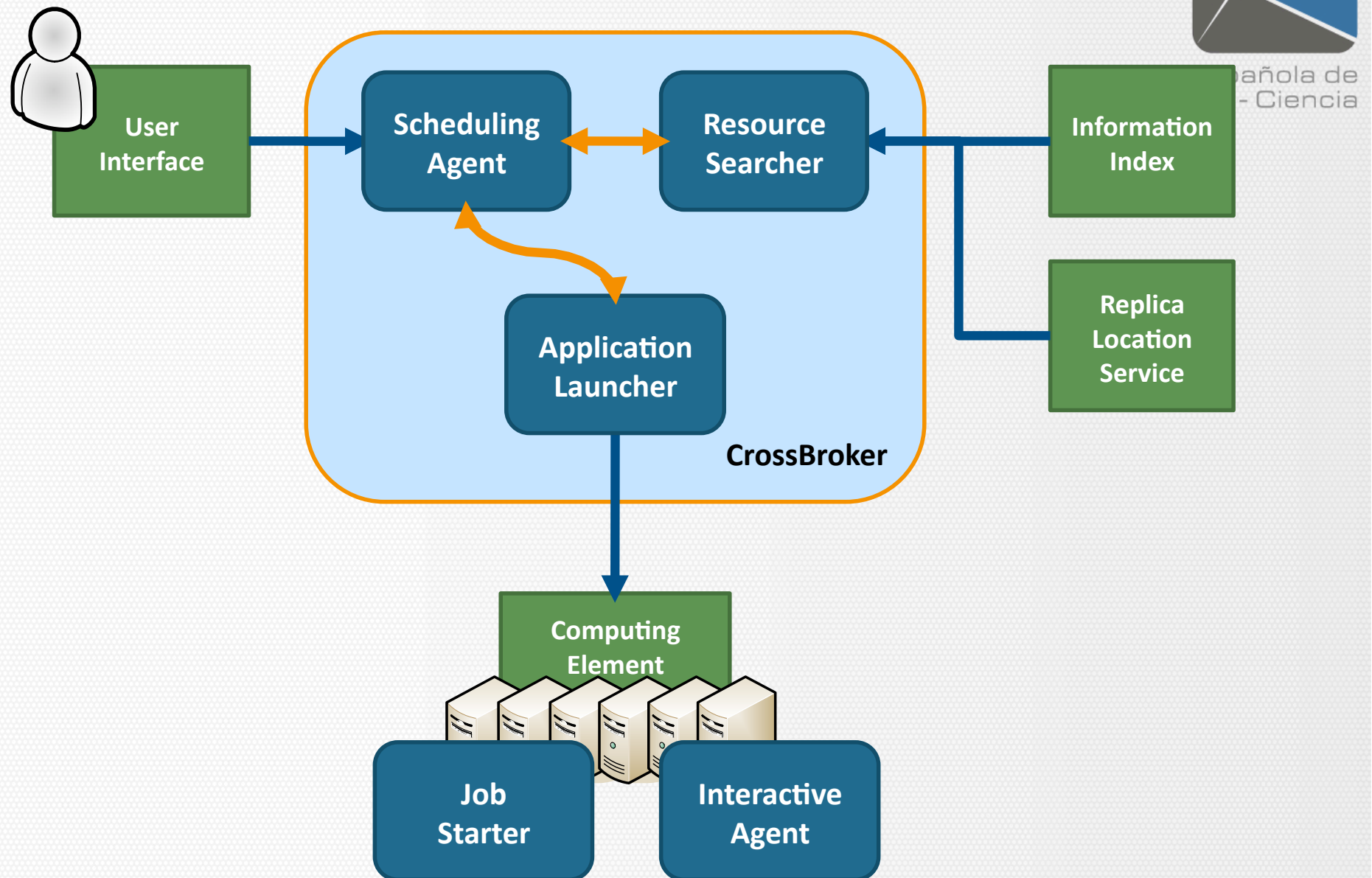- More information at: http://www.gridway.org

# Metaschedulers: GridWay

# Metaschedulers: CrossBroker

- CrossBroker provides support for Parallel and Interactive jobs
  - High priority treatment for interactive jobs with the use of multiprogramming
  - Interoperable with EGEE, provides same services than gLite WMS

- Used in production environments
  - Used in EU CrossGrid, int.eu.grid and Euforia projects (12K – 55K jobs per month)

- Developed by UAB + CSIC
- More info: http://www.oliba.uab.es/crossbroker

# Metaschedulers: CrossBroker

# Interactivity Support

- Interactivity allows researchers to visualize results and obtain them faster

- Requirements:
  - Fast startup: the possibility of starting the application immediately, even in high occupancy scenarios
  - Online Input-Output streaming: the ability to have application input and output online.
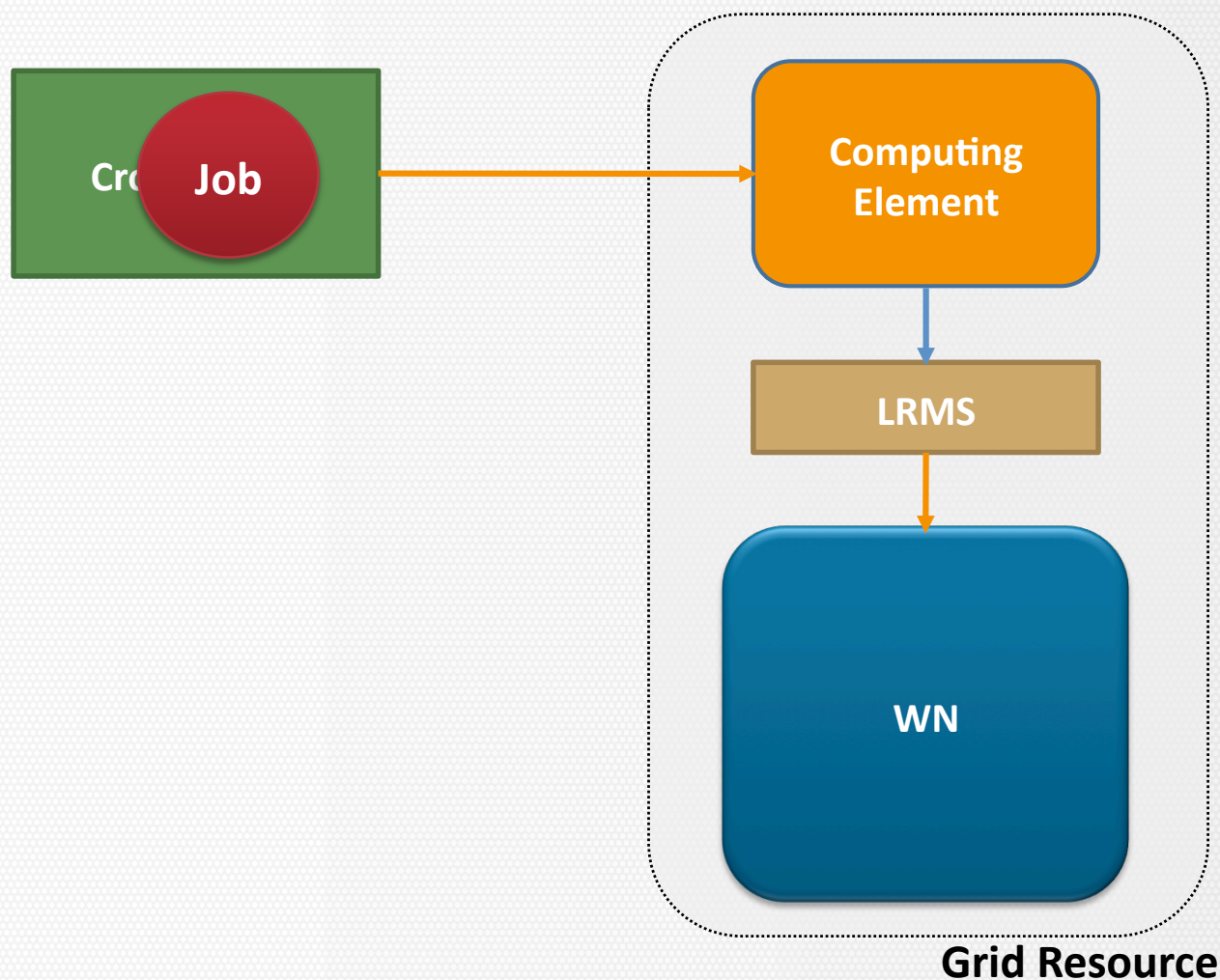
# Interactivity Support: Multiprogramming

- The idea
  - Each job is encapsulated in an agent that takes control over the WN independently of its LRMS

- Lightweight "Virtual Machines"
  - Each Worker Node is divided in 2 execution slots
  - Each VM can execute jobs independently (e.g. batch and interactive)
  - NOT a full virtual machine (Xen, VMWare,…)
  - NO need for special priviledges in the WN

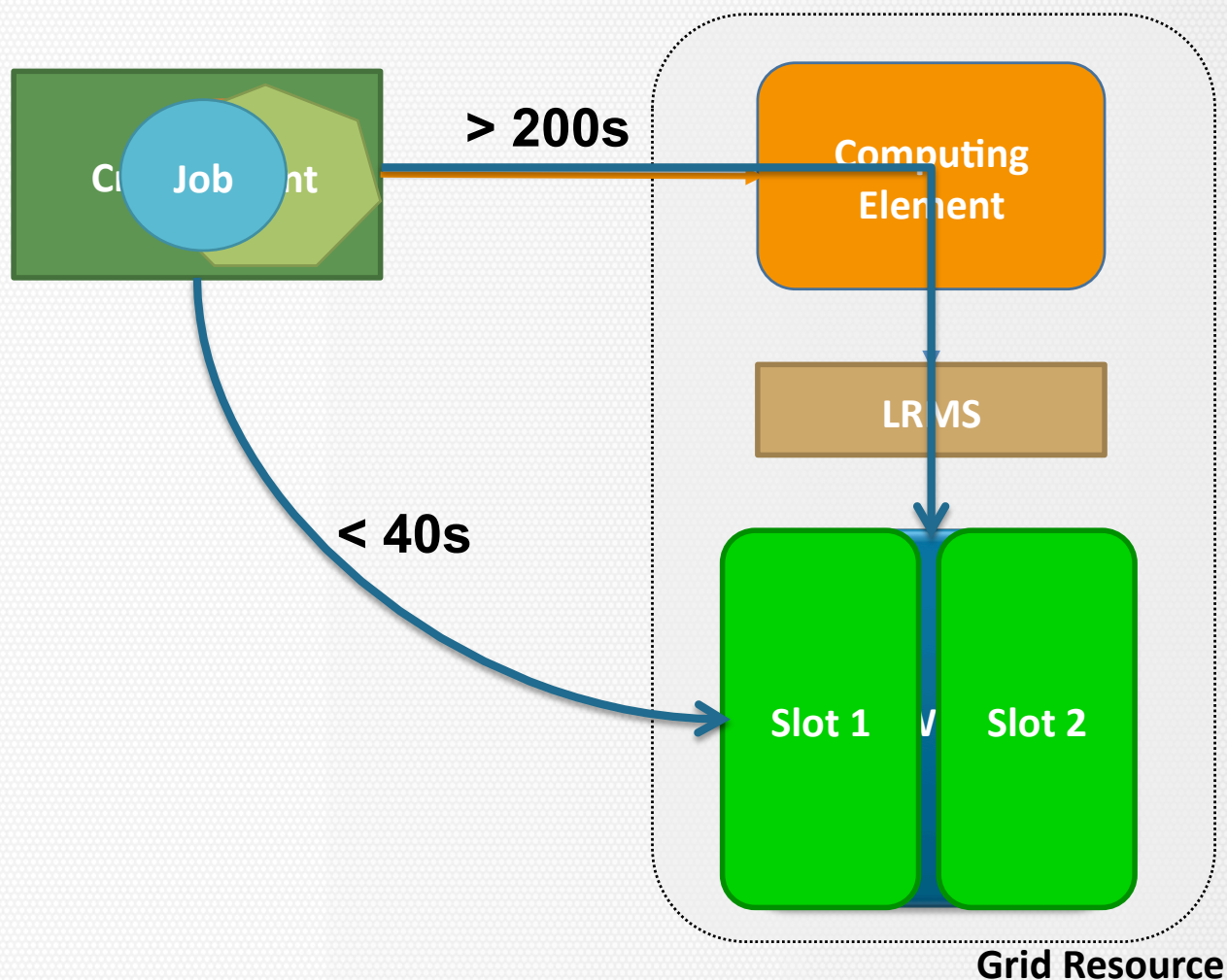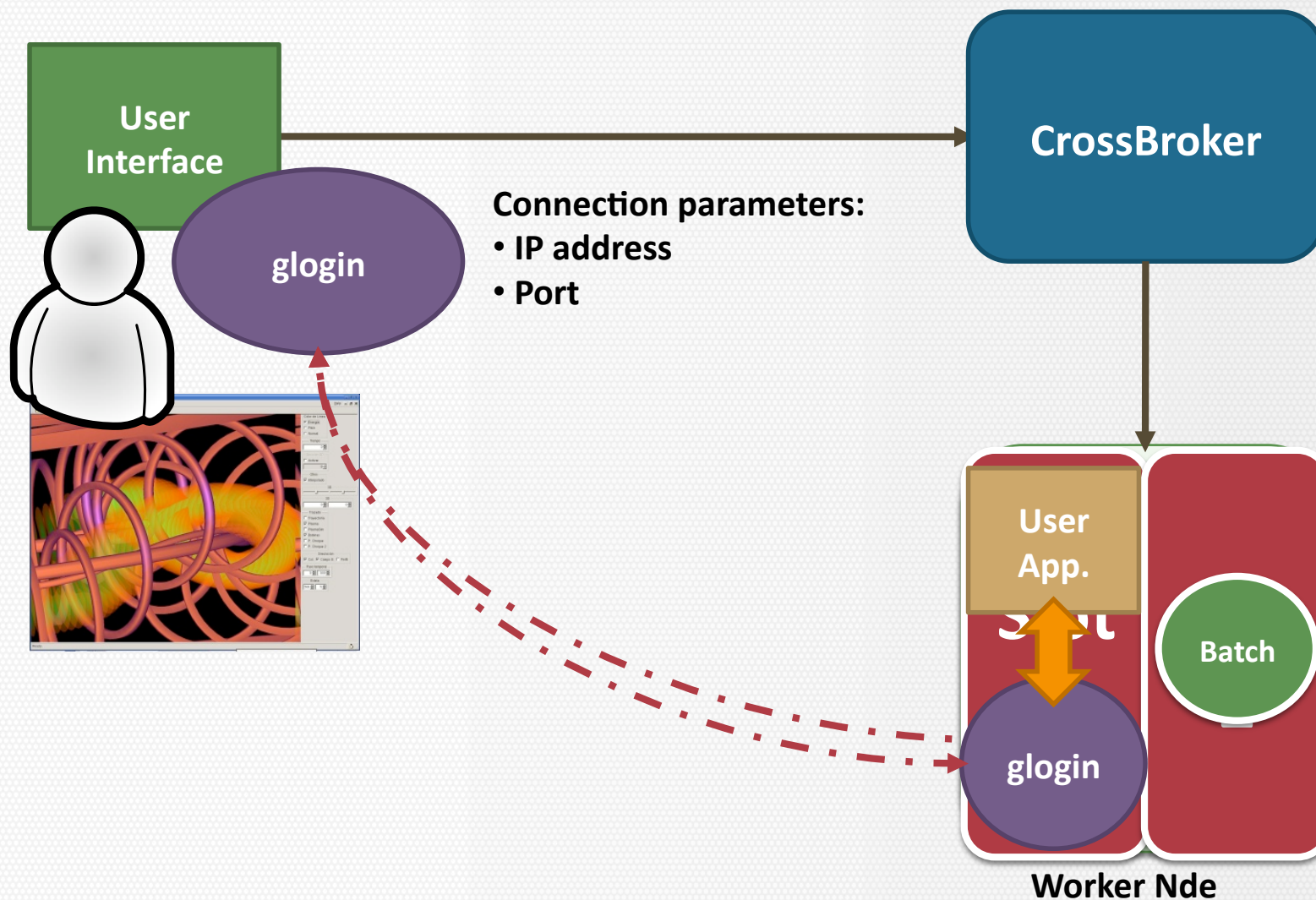# Interactivity Support: Multiprogramming

# Interactivity Support: Multiprogramming

# Interactivity Support: Interactive Agents

**User Interface**

**glogin**

Connection parameters:
- **IP address**
- **Port**

**CrossBroker**

**User App.**

**Batch**

**glogin**

**Worker Nde**

# MPI Support

- Many application areas need MPI support
  - Earth Sciences, Biological sciences, Computational Chemistry, Nuclear Fusion,
  - Representative results can be obtained by using order of 10s-100s of CPUs
- Many clusters are MPI – ready
  - In local mode by direct submission
  - Shared filesystems with high performance intranet
- It is interesting to offer this capability when the user is working inside a Grid infrastructure
  - As an infrastructure on its own
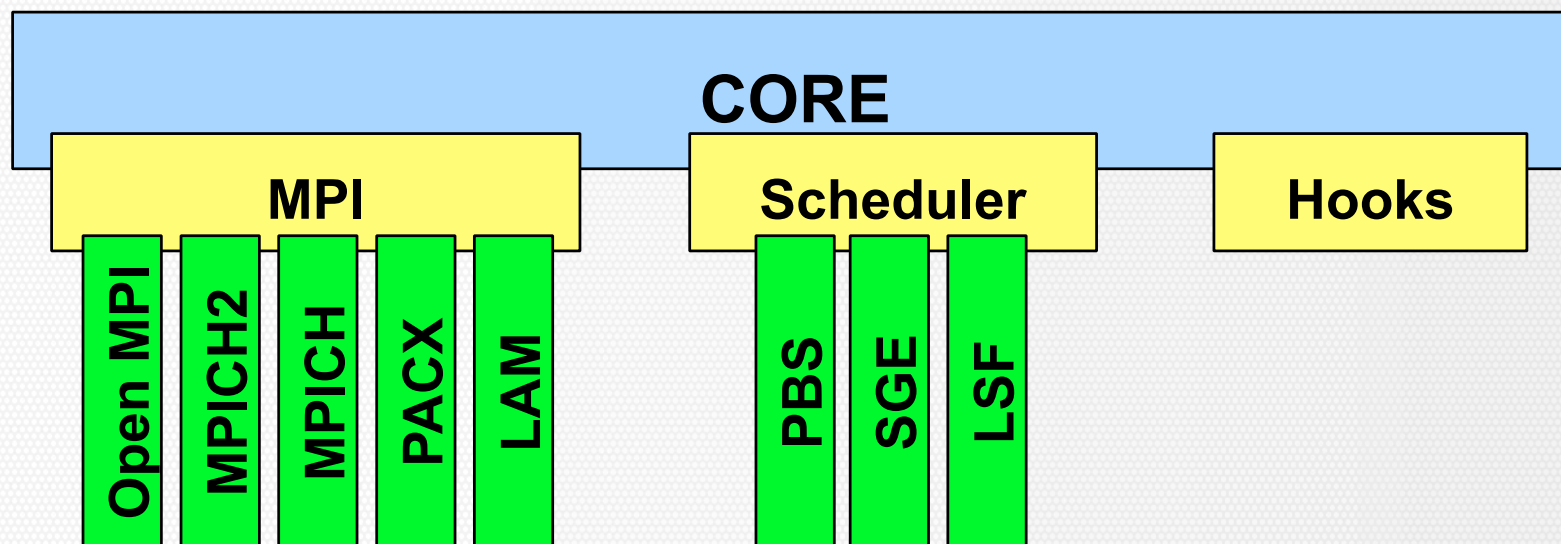  - As a testbed for small runs before executing on large HPC

# MPI Support: Issues

- There is no standard way of starting an MPI app
  - No common syntax for mpirun, mpiexec support optional
  - Schedulers (PBS, SGE, Condor…) handle machinefiles in different ways
  - Metascheduler services need to handle different implementations (OpenMPI, MPICH, LAM…) in a flexible and extensible way
  - Shared vs. Non-Shared filesystems

# MPI Support: MPI-Start

- Specifies a unique interface to the upper layer in the middleware to describe MPI jobs

- Support basic file distributions

- Implemented as portable shell scripts

- Extensible via user hooks and plugins at the site level

**CORE**

**MPI** — Open MPI | MPICH2 | MPICH | PACX | LAM

**Scheduler** — PBS | SGE | LSF

**Hooks**

# MPI Support: MPI-Start

- MPI-Start is used by the CrossBroker to support:
  - Intra-cluster apps with OpenMPI or MPICH
  - Inter-cluster apps using PACX-MPI or MPICH-G2
- User does not need to specify low level details of jobs:
  - MPI-Start copies all the input files to the WN (if not shared FS)
  - MPI-Start invokes the mpi program using the best configuration for the site.

# Summary

- NGI Grid Infrastructure is well active!
  - Deployment of NGI infrastructure going on
  - Integrates EGEE like resources (i2g, EUFORIA, DORII, EELA, …) and GT 4 resources
- Added values of NGI-ES
  - Metaschedulers middleware
  - Support to MPI and Interactive jobs
- Interest in European Grid Infrastructure
- More info: www.e-ciencia.es