# Grid Job Submission Concepts And Basic Examples

**Daniele Cesini**
**Marco Bencivenni**
**INFN-CNAF**
**Grid School**
**Bologna, 21-24 February 2011**

# Overview

- **General Grid Job Submission concepts**
  - What is the Grid from a job submission perspective
  - Involved Actors
    - UI, WMS, LB
  - Submission models
  - Input and Output
  - Resource Selection
  - Grid Failures/Error Recovery
- **The Grid Job**
  - Anatomy of a Grid Job
  - The Job State Machine
  - The Grid JOB ID
  - Job Types
- **Hands On**
  - Let's submit our first Grid job…

**Ian Foster's "What is the Grid? A Three Point Checklist" [3]**

1) coordinates resources that are not subject to centralized control
2) *using standard, open, general-purpose protocols and interfaces*
3) to deliver nontrivial qualities of service
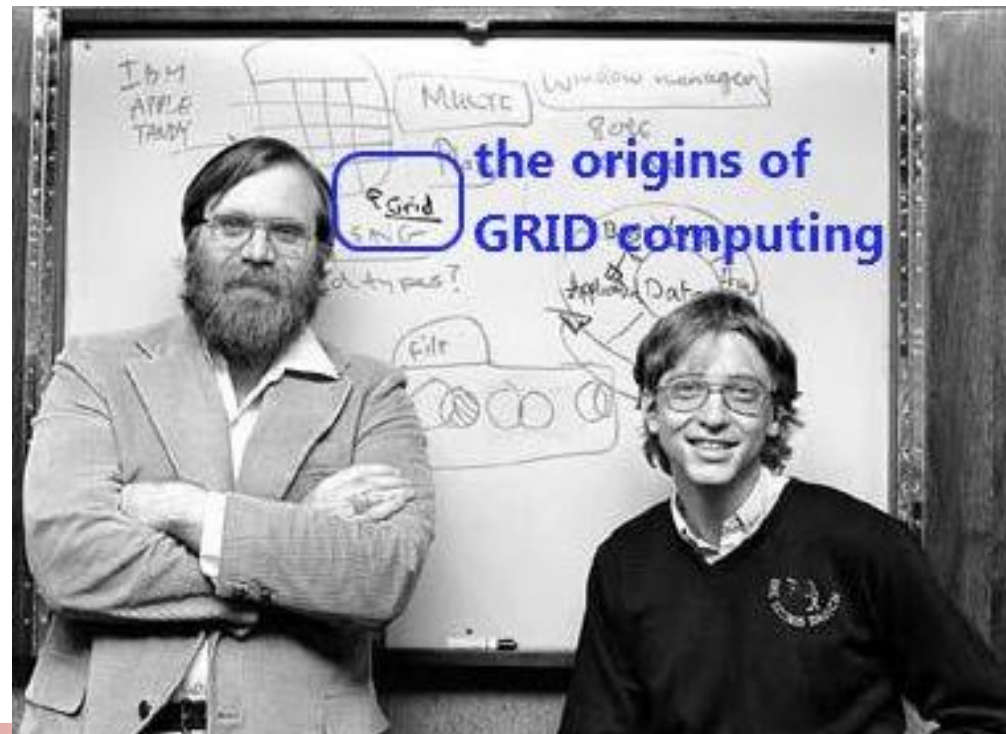
Ian Foster

Carl Kesselman

[1] Foster, I. and Kesselman, C. eds. The Grid: Blueprint for a New Computing Infrastructure, Morgan Kaufmann, 1999, 259-278

[2] Ian Foster, Carl Kesselman, and Steven Tuecke. 2001. The Anatomy of the Grid: Enabling Scalable Virtual Organizations. Int. J. High Perform. Comput. Appl. 15, 3 (August 2001), 200-222. DOI=10.1177/109434200101500302
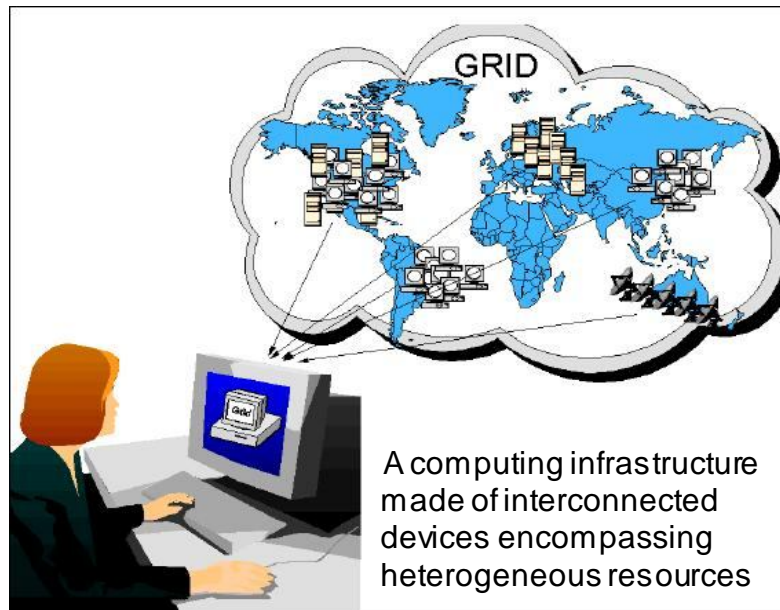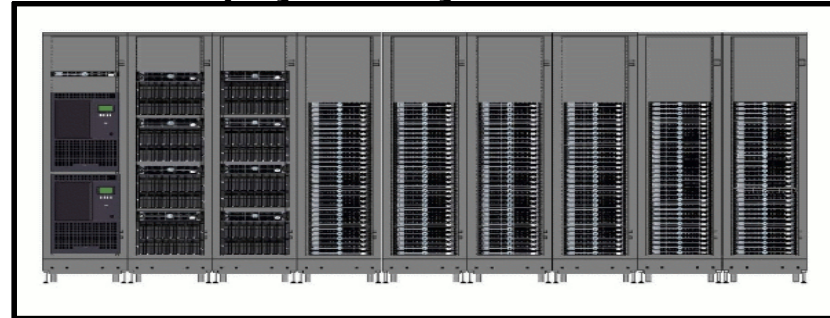
[3] What is the Grid? A Three Point Checklist. I. Foster, GRIDToday, July 20, 2002.

# 1) No centralised control

The user in general has full ownership of a desktop workstation.

A Cluster is a shared resource – Only the administrator has full control of the system
The physical layer is well defined

A computing infrastructure made of interconnected devices encompassing heterogeneous resources

In a Grid both users and physical layer are (should be) virtualised

I submit my jobs to "the GRID" and they get processed: somehow, somewhere, after some time.

The is no GRID owner!
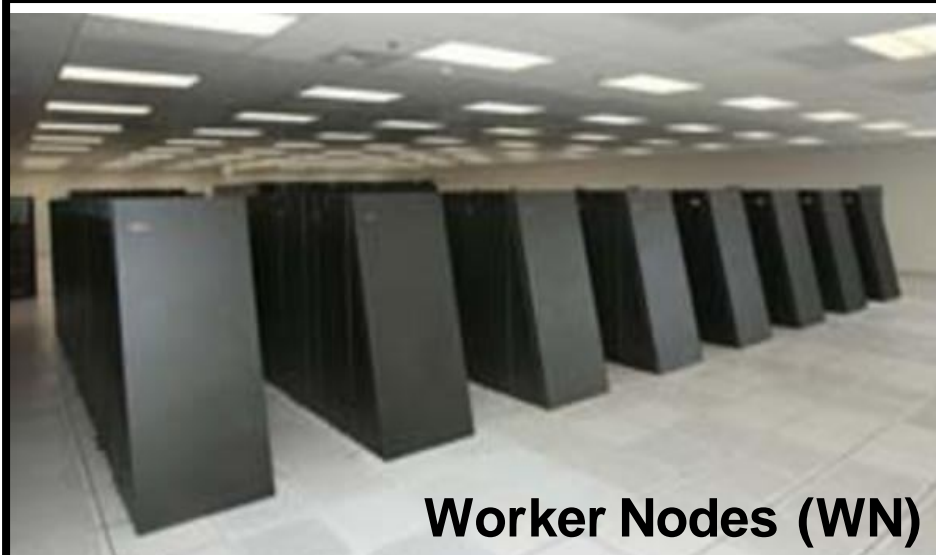
**In this school Grid means the WLCG/EGI gLite based Grid**

**Job submission is done through the gLite WMS or via direct access to the CREAM Computing Element**

gLite Home: http://www.glite.org
WMS Project: http://web.infn.it/gLiteWMS/
CREAM Project: http://grid.pd.infn.it/cream/
EGI Project: http://www.egi.eu

**Worker Nodes (WN)**

**Storage**

**Batch System (PBS,LSF…)**

**Computer Element (CE)**

**Storage Element (SE)**

**Site BDII (sBDII)**

PRODUCTION Physical and Logical CPUs Per Date
January-2004 to February-2011

*Jan 2005 – Feb 2011*

(Source EGI accounting portal)

**Tianhe-1A:**
**14,336 Xeon X5670**
**7,168   Tesla M2050**

**2566 Tflops max**
**4701 Tflops peak**

**Bull Tera-100:**
**140,000 Intel Xeon**

**1050 Tflops max**
**1255 Tflops peak**

**HPC**

**(Source top500.org)**

**SETI  HOME**

$$\sim 3 \cdot 10^5 \cdot 3 \cdot 10^9 \sim 9 \cdot 10^{14} \sim 900 \text{Tflops}$$

**HTC = High throughput computing**

**2,829,110 Hosts**
**486 Tflops**
**(source boincstats.com 02/2011)**

# The Job Submission Actors

UI

High Level Application

Information Systems (IS)

WMS

Logging & Bookkeeping (LB)

VOMS

MYPROXY

LCG File Catalog (LFC)

# The Job Submission Actors

**Grid Resources:** …well, it's the Grid…

**User:** …well you know who you are…

**UI:** a machine containing a collection of clients to access the Grid services

**WMS:** responsible for dispatching user jobs selecting the best possible resource according to the job requirements

**LB:** contains detailed information about jobs lifecycle – thigtly coupled to the WMS

**IS:** information system – contains an updated snapshot of what is contained in the Grid

**LFC:** a file catalog, links logical file names to physical locations

**VOMS & Myproxy:** security stuff to obtain valid certificates

# Overall architecture



© Maarten Litmaath, CERN

# Scheduling models

Scheduling of distributed, data-driven applications in a Grid environment is a challenging problem. From the initial design two submitting models have been foreseen:

## eager scheduling ("push" model)

- The job is bound to a resource as soon as possible.
- Once the decision has been taken, the job is passed to the selected resource for execution
- It will probably end up in a queue.

## lazy scheduling ("pull" model)

- The job is held by the WMS until a resource becomes available.
- When this happens the resource is matched against the submitted job.

Currently only the push mode is adopted.
LHC VO are moving towards "pull models" built on top of the native "push implementation"

- **Direct submission to computing resources**
  - WMS bypassed
  - Users cannot have a global view of the whole picture
  - The responsibility of the job remains to the user, most of the time becoming a 'burden'

- **Submission through the WMS**
  - passes the responsibility of the outcome of the job to this service which provides value-added capabilities and instrumented to always know the whole picture (i.e. Non-trivial QoS).
  - Match Making can be avoided  using the –r option
    - Other capabilities preserved

# ISB & OSB through WMS



**OUTPUT SANDBOX (OSB)**

**INPUT SANDBOX (ISB)**

YAY! IT WORKS!!!

**ISB & OSB**

**Can travel with the job through the WMS**
➔ **MUST BE SMALL (~5MB) !!!**

**Can be remote** ➔ **Any size**
**USE DATA MANAGEMENT COMMANDS**

# WMS Supported Job Types

- **Batch-like**

- **DAG workflow**

- **Collection**

- **Parametric**

  k

- **MPI**

compound

La Gare Montparnasse, 1895

"A distributed system is one in which the failure of a computer you didn't even know existed can render your own computer unusable."
Leslie Lamport

# Expect the unexpected

- **When services / servers don't respond or return an invalid status / message;**

- **When the air-conditioning / power fails (again & again & again);**

- **When disks fail and you have to recover from backup – and the tapes have been overwritten;**

- **When a service engineer puts a Coke into a machine to 'warm it up'**

- **When Oracle returns you someone's else data**

- **When a fishing trawler cuts a trans-Atlantic network cable;**

- **When a Tsunami does the equivalent in Asia Pacific;**

**All these things really happened   ©Jamie Shiers**
**2008 J. Phys.: Conf. Ser. 119 052030**

**The Grid (i.e. the WMS) can recover Grid errors**

**Executable errors are not Grid errors!!**

## Errors are recovered through resubmission

### Deep resubmission:

- When the user's job fails after having started running on the WN
- On every grid failure (even before the job started on the WN) if the shallow is disabled
- May be problematic if the job touches data

### Shallow resubmission:

- If failed before having started the execution on the WN
- Safer than the deep

# WMS Summary

- **Implements a "push" submission model**

- **Hides the Grid complexity to the user**

- **Takes charge of completing the job**

- **Select (hopefully) the best resource for the job**
  - Based on user defined criteria

- **Offers Sandbox management capabilities**

- **Implements error recovery capabilities**

- **Allows to submit multiple job types and workflows**

# The single batch Grid Job

| | |
|---|---|
| **JOB Type** | JobType = "Normal"; |
| **Prologue** | Prologue = "prologue.sh"; |
| **Input SandBox** | InputSandbox = {"test.sh", "fileA"}; |
| **Requirements** | Requirements = false; |
| **Executable** | Executable = "test.sh"; |
| **Std Output/Error** | StdOutput = "std.out";<br>StdError = "std.err"; |
| **Output SandBox** | OutputSandbox={"std.out", "std.err"}; |
| **Epilogue** | Epilogue = "compress.sh"; |
| **Error Recovery** | RetryCount = 1;<br>ShallowRetryCount = 2; |

**JDL for compound jobs are based on normal job JDL**

# The GRID Job Identity Card

```
==================== glite-wms-job-submit Success ====================

The job has been successfully submitted to the WMProxy
Your job identifier is:

https://lb-server-03.cnaf.infn.it:9000/C-Et5jbMMBjjUHkT1X6wVg

==================================================================
```

## JobID:

- Upon submission each job is assigned a unique, virtually non-recyclable job identifier in an URL form.

- The server part of the URL designates the **LB server**

- The remainder is a random generated sequence: the Grid is a highly decentralized system, characterized by lack of unified control ➔ **no serial numbering is possible**

**Submitted** The job has been submitted from the UI but it is still waiting to be accepted by the WMS

**Waiting** The job has been accepted by the WMS and it is waiting to be processed fro the Match Making

**Ready** The job has been processed by the WMS but it hasn't been transferred to the CE yet

Scheduled job is waiting in the CE queue

Running job is executing!

**Done**  The job has terminated, either successfully or considered to be terminated with some error. (i. e.: due to unrecoverable errors on the CE side)

**Aborted** The processed job has been aborted by the WMS

(for too long in a queue on the WM or on the CE, expired credentials etc.)

**Cancelled** Job has been cancelled by the user

**Cleared:** the output has been transferred by the user or removed because of some timeout

# First of all.....the proxy!

**Make sure you have your certificate in the .globus dir**

```
[cesini@ui cesini]$ ll .globus/

-rw-------    1 cesini   cesini      2126 Jul  7  2007 usercert.pem
-r--------    1 cesini   cesini      1910 Jul  7  2007 userkey.pem
```

## voms-proxy-init --voms <vo_name>

## voms-proxy-info --all

> [cesini@ui corso]$ cat minimal.jdl
> Executable = "/bin/hostname";
> StdOutput = "std.out";
> StdError = "std.err";

**HANDS ON**

**Submit this minimal, completely useless JDL**

**Why is it useless??**

> *glite-wms-job-submit -a minimal.jdl*

**If it's not working don't worry...**

# OK submission output

[cesini@ui corso]$ glite-wms-job-submit -a minimal.jdl

Connecting to the service https://glite-rb-00.cnaf.infn.it:7443/glite_wms_wmproxy_server

=================glite-wms-job-submit Success ============
The job has been successfully submitted to the WMProxy
Your job identifier is:

https://lb009.cnaf.infn.it:9000/TWr2bZ0QIaWsBrd43zslAg

===================================================

**UNIQUE JOB ID**
**LB host in the JOBID**

**-e  <server endpoint>** allows to override the default server used by the client

**-c <conf_file_name>** allows to use a custom config file

**Exercise: try to change the WMS endpoint with both options**

**Hint1: you first need to discover which are the WMS available to your VO (lcg-infosites –vo <your_vo>  wms)**

**Hint2: a config file is on the course UI in the "corso" folder**

# The Job Status

*[cesini@ui corso]$* **glite-wms-job-status**
   *https://albalonga.cnaf.infn.it:9000/TWr2bZ0QIaWsBrd43zslAg*

***************************************************************

*BOOKKEEPING INFORMATION:*

*Status info for the Job :*
   *https://albalonga.cnaf.infn.it:9000/TWr2bZ0QIaWsBrd43zslAg*
*Current Status:      Ready*
*Destination:        grid003.roma2.infn.it:2119/jobmanager-lcgpbs-cert*
*Submitted:          Mon Nov 19 15:09:42 2007 CET*

***************************************************************

**Try to increase the output verbosity**
**Hint: Use –v <1|2|3>**

**Try to open the JobID URL!!!**

**Exercise: try to submit a JDL with ISB and OSB**

**Hint: Locate the first.jdl file in your UI**

```
[cesini@ui corso]$ cat first.jdl
##################################################
# My First JDL with very basic attributes #
##################################################

Executable = "test.sh";
Arguments = "fileA fileB";
StdOutput = "std.out";
StdError = "std.err";
InputSandbox = {"test.sh", "fileA", "fileB"};
OutputSandbox = {"std.out", "std.err"};
```

# "-i" and "-o"

**Exercise: try to submit many jobs saving the JobID to be used later on**

**Hint: use "-o" when submitting and "-i" with the status command**

```
[cesini@ui corso]$ export i=0
[cesini@ui corso]$ while [ $i -le 10 ]; do glite-wms-job-submit –a –c
wms_rb00.conf -o ID_file.txt first.jdl ; let i=i+1; done >> submission.txt &
```

```
[cesini@lcg-ui corso]$ glite-wms-job-status -i ID_file.txt
------------------------------------------------------------------
1 : https://albalonga.cnaf.infn.it:9000/8FjA0EJ05jYHdkgYX0JU3Q
...........
10: https://albalonga.cnaf.infn.it:9000/5ZWpn7uomUzXtjqxFxJe5g
11: https://albalonga.cnaf.infn.it:9000/kdaNgNOSEwHzlzV47K1Fwg
a : all
q : quit
------------------------------------------------------------------
Choose one or more jobId(s) in the list - [1-11]all:
```

**Use - -noint to have directly the status for all the JOBIDs**

# What is that "-a" ?

**-a** means automatic delegation of the proxy to the WMS
- is handy
- is SLOW – each job submission implies an SSL delegation

**-d <name>** means use a pre delegated proxy
- you need to pre delegate a proxy with name <name>
- is FASTER

glite-wms-job-delegate-proxy --help

Exercise: try to test the submission timing with both options

Hint: Solution in the next slide

# Proxy delegation

```
[cesini@gridlab20 corso]$ glite-wms-job-delegate-proxy -d pippo
Connecting to the service https://prod-wms-01.pd.infn.it:7443/glite_wms_wmproxy_server
================== glite-wms-job-delegate-proxy Success ==================
Your proxy has been successfully delegated to the WMProxy(s):
https://prod-wms-01.pd.infn.it:7443/glite_wms_wmproxy_server
with the delegation identifier:  pippo

=========================================================================
```

```
[cesini@gridlab20 corso]$ time glite-wms-job-submit -a -e https://prod-wms-
     01.pd.infn.it:7443/glite_wms_wmproxy_server minimal.jdl
Connecting to the service https://prod-wms-01.pd.infn.it:7443/glite_wms_wmproxy_server
==================== glite-wms-job-submit Success ====================
The job has been successfully submitted to the WMProxy
Your job identifier is:
https://prod-lb-01.pd.infn.it:9000/XXYeKKsiwWGfWPfgVxBZVw

=========================================================================
```

**real   0m1.027s**
user    0m0.151s
sys     0m0.013s

```
[cesini@gridlab20 corso]$ time glite-wms-job-submit -d pippo -e https://prod-wms-
     01.pd.infn.it:7443/glite_wms_wmproxy_server minimal.jdl
Connecting to the service https://prod-wms-01.pd.infn.it:7443/glite_wms_wmproxy_server
==================== glite-wms-job-submit Success ====================
The job has been successfully submitted to the WMProxy
Your job identifier is:
https://prod-lb-01.pd.infn.it:9000/0ux6eOR-Kanrr0wV-anXRg

=========================================================================
```

**real   0m0.735s**
user    0m0.090s
sys     0m0.016s

# Now I need the output!

**[cesini@ui corso]$ glite-wms-job-output**
**https://albalonga.cnaf.infn.it:9000/TWr2bZ0QIaWsBrd43zslAg**

**Connecting to the service**
**https://131.154.100.90:7443/glite_wms_wmproxy_server**
==================================================================
**JOB GET OUTPUT OUTCOME**
**Output sandbox files for the job:**
**https://albalonga.cnaf.infn.it:9000/TWr2bZ0QIaWsBrd43zslAg**
**have been successfully retrieved and stored in the** directory:
**/tmp/glite/glite-ui/cesini_TWr2bZ0QIaWsBrd43zslAg**
==================================================================

**Exercise1: try to retrieve the output of:**
**1) the first.jdl done job**
**2) the first.jdl job not done yet**
**3) the minimal.jdl done job**

**Exercise2: change the default target dir of the output**
**Hint: " - -dir"**

# I need to cancel my job!

*[cesini @ui corso]$ **glite-wms-job-cancel***

*https://albalonga.cnaf.infn.it:9000/kFTSkFWGadkZqFgNb4m5WQ*

*Are you sure you want to remove specified job(s) [y/n]y : y*

*Connecting to the service
    https://131.154.100.90:7443/glite_wms_wmproxy_server*

*============== glite-wms-job-cancel Success ===================*

*The cancellation request has been successfully submitted for the following job(s):*

*- https://albalonga.cnaf.infn.it:9000/kFTSkFWGadkZqFgNb4m5WQ*

*===========================================================*

**Exercise: try to cancel one of your job:**
 1) **Before it's done**
 2) **When it's done**

```
[cesini@ui corso]$ glite-wms-job-logging-info https://albalonga.cnaf.infn.it:9000/fzxo1li1K-
            sCjAFfHIjQ3Q
**********************************************************
LOGGING INFORMATION:
Printing info for the Job : https://albalonga.cnaf.infn.it:9000/fzxo1li1K-sCjAFfHIjQ3Q
     ---
Event: RegJob
- source               =   NetworkServer
- timestamp            =   Mon Nov 19 15:27:36 2007 CET
     ---
Event: RegJob
- source               =   NetworkServer
- timestamp            =   Mon Nov 19 15:27:36 2007 CET
     ---
Event: Accepted
- source               =   NetworkServer
- timestamp            =   Mon Nov 19 15:27:37 2007 CET
     ---
Event: EnQueued
- result        =   START
- source        =   NetworkServer
- timestamp     =   Mon Nov 19 15:27:37 2007 CET
     ---
Event: EnQueued
- result        =   OK
- source        =   NetworkServer
- timestamp     =   Mon Nov 19 15:27:37 2007 CET
     ---
Event: DeQueued
- source               =   WorkloadManager
- timestamp            =   Mon Nov 19 15:27:37 2007 CET
     ---
Event: Match
- dest_id       =   spacin-ce1.dma.unina.it:2119/jobmanager-lcgpbs-cert
- source        =   WorkloadManager
- timestamp     =   Mon Nov 19 15:27:41 2007 CET
```

**Try to increase verbosity up to –v 3**

**Identify the JDL in the output and compare with your JDL**

# References

- Ian Foster, Carl Kesselman, and Steven Tuecke. 2001. The Anatomy of the Grid: Enabling Scalable Virtual Organizations. *Int. J. High Perform. Comput. Appl.* 15, 3 (August 2001), 200-222. DOI=10.1177/109434200101500302

- What is the Grid? A Three Point Checklist. I. Foster, GRIDToday, July 20, 2002.

- The Grid: A New Infrastructure for 21st Century Science. I. Foster. *Physics Today*, 55(2):42-47, 2002.

- The Physiology of the Grid: An Open Grid Services Architecture for Distributed Systems Integration. I. Foster, C. Kesselman, J. Nick, S. Tuecke, Open Grid Service Infrastructure WG, Global Grid Forum, June 22, 2002.

# Useful Links

- **WMS Project Homepage**
  http://web.infn.it/gLiteWMS/

- **WMProxy submission**
  https://edms.cern.ch/document/590869/1

- **LB documentation**
  http://egee.cesnet.cz/en/JRA1/LB/documentation.php

- **Glite UserGuide**
  https://edms.cern.ch/file/722398//gLite-3-UserGuide.html

- **Glite General Documentation Page**
  http://glite.web.cern.ch/glite/documentation/

# Useful Links

- **CREAM Project home page**

  http://grid.pd.infn.it/cream/

- **Investigating Job Submission Description Language (JSDL)**

  https://forge.gridforum.org/projects/jsdl-wg/

- **Condor ClassAdd**

  http://www.cs.wisc.edu/condor/classad/refman/

- **MPI in gLite**

  http://grid.ie/mpi/wiki/JobSubmission

  http://egee-uig.web.cern.ch/egee-uig/production_pages/MPIJobs.html

# Break!