



Enabling Grids for E-scienceE

# Advanced Job Submission on the Grid

Antun Balaz

Scientific Computing Laboratory

Institute of Physics Belgrade

<http://www.scl.rs/>



30 Nov – 11 Dec 2009

[www.eu-egee.org](http://www.eu-egee.org)

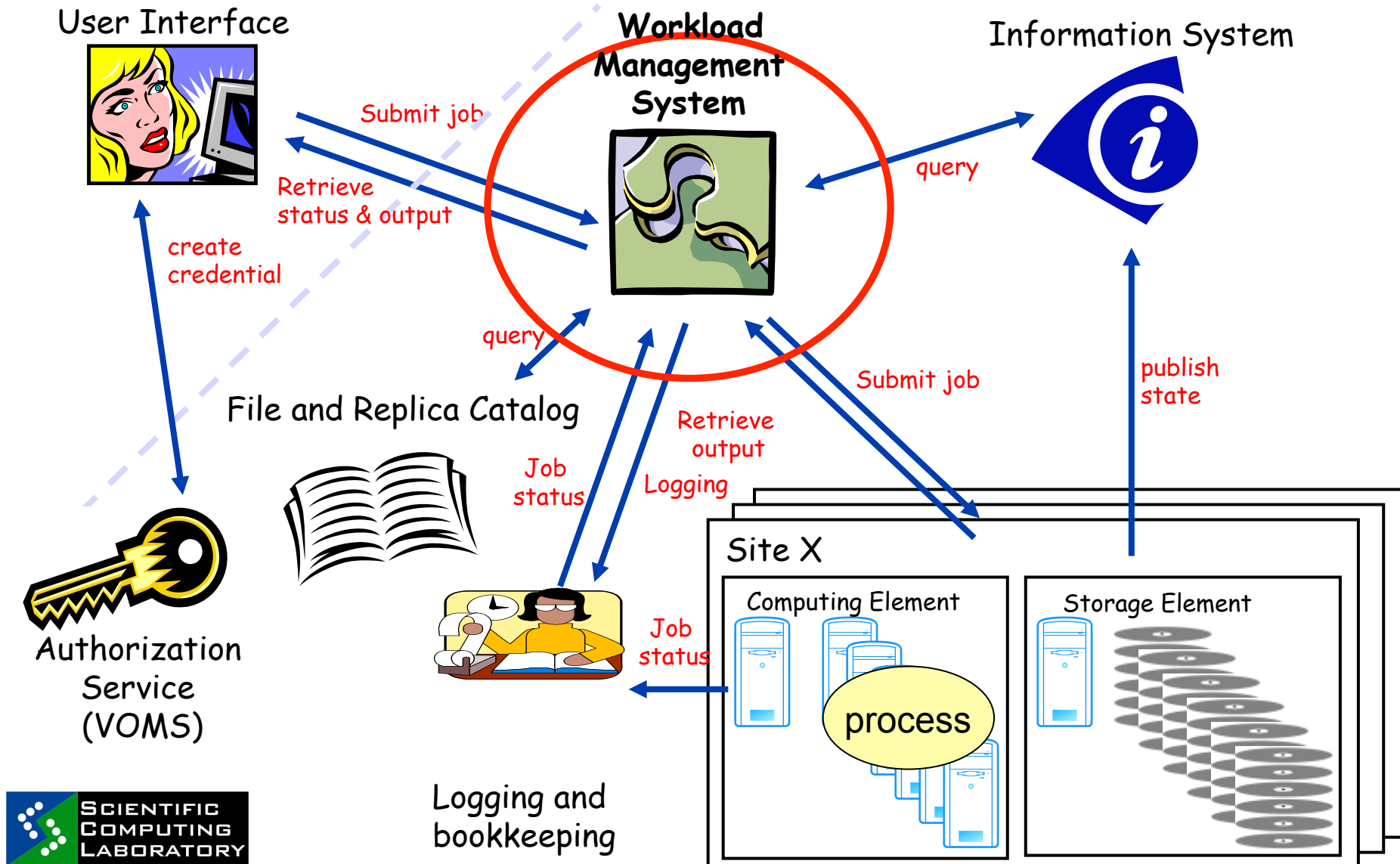


SEE-GRID-SCI  
SEE-GRID eInfrastructure for regional eScience



Information Society

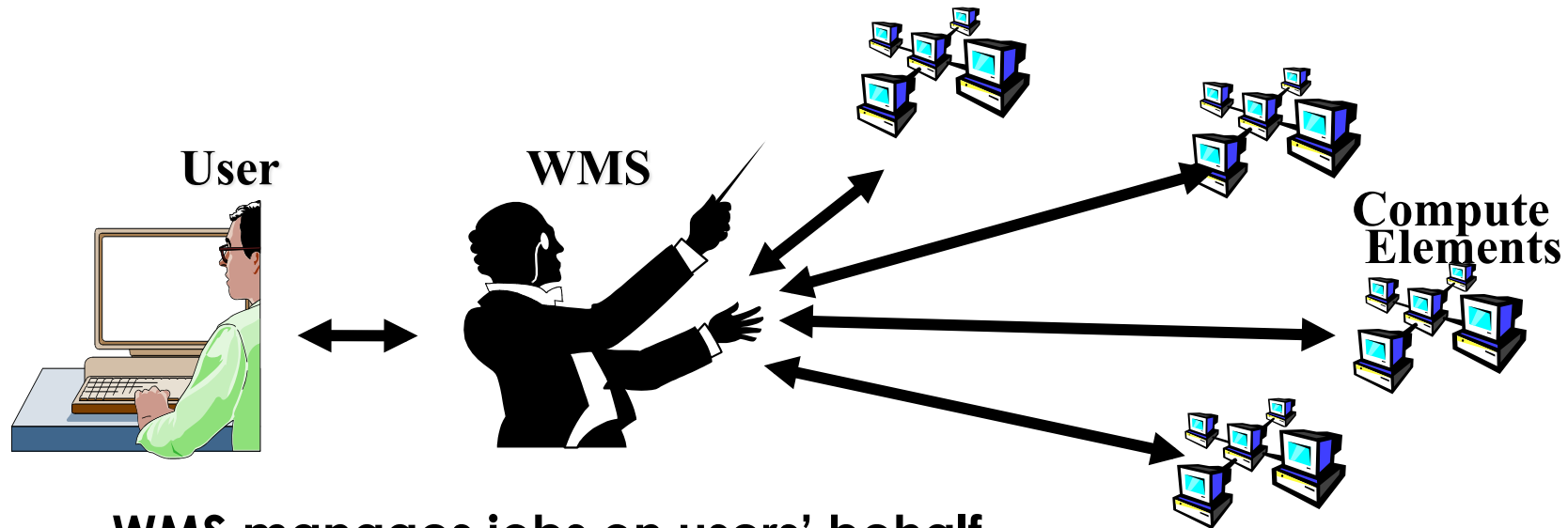




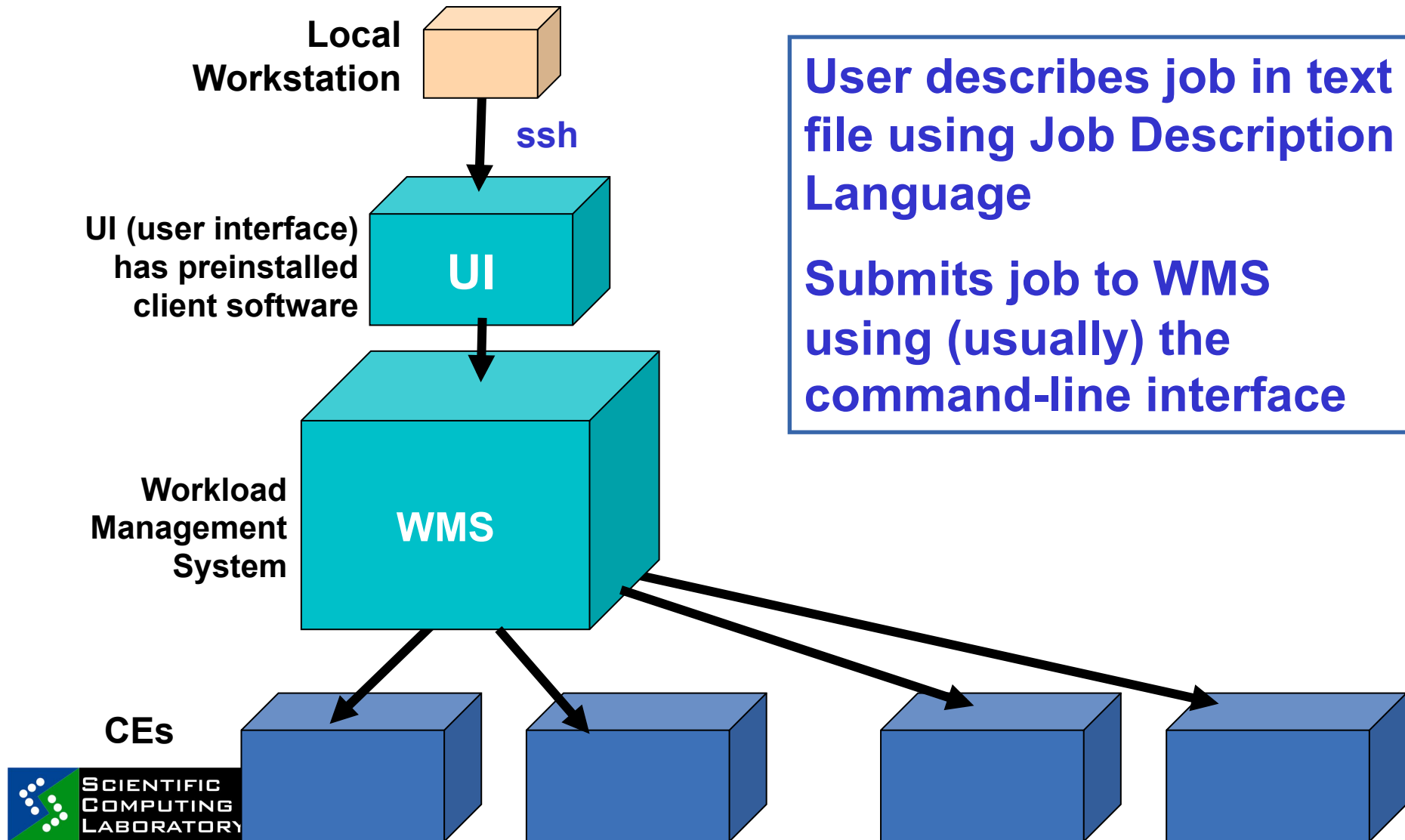
- **Again a little bit on WMS**
- **How are your jobs handled**
- **JDL attributes**
- **Which types of jobs exist**
- **Examples of complex JDLs**

## Why does the Workload Management System exist?

- **Grids have**
  - Many users
  - Many jobs – a “job” = an executable you want to run
  - Where many compute nodes are available
  - Workload Management System is a software service that makes running jobs easier for the user
- **It builds on the basic grid services**
  - E.g. Authorisation, Authentication, Security, Information Systems, Job submission
- **Terminology: “Compute element”:** defined as a batch queue - One cluster can have many queues



- **WMS manages jobs on users' behalf**
  - User doesn't decide where jobs are run
  - User defines the job and its requirements, WMS matches this with available CEs
- **Effect:**
  - Easier submission
  - Users insulated from change in Compute elements
  - WMS – can optimise your jobs – e.g. which CE?



- Jobs run in batch mode on grids.
- Steps in running a job on a gLite grid with WMS:
  1. Create a text file in “Job Description Language”
  2. Optional check: list the compute elements that match your requirements (“list match” command)
  3. Submit the job ~ “glite-wms-job-submit -a myfile.jdl”  
Non-blocking - Each job is given an id.
  4. Occasionally check the status of your job
  5. When “Done” retrieve output

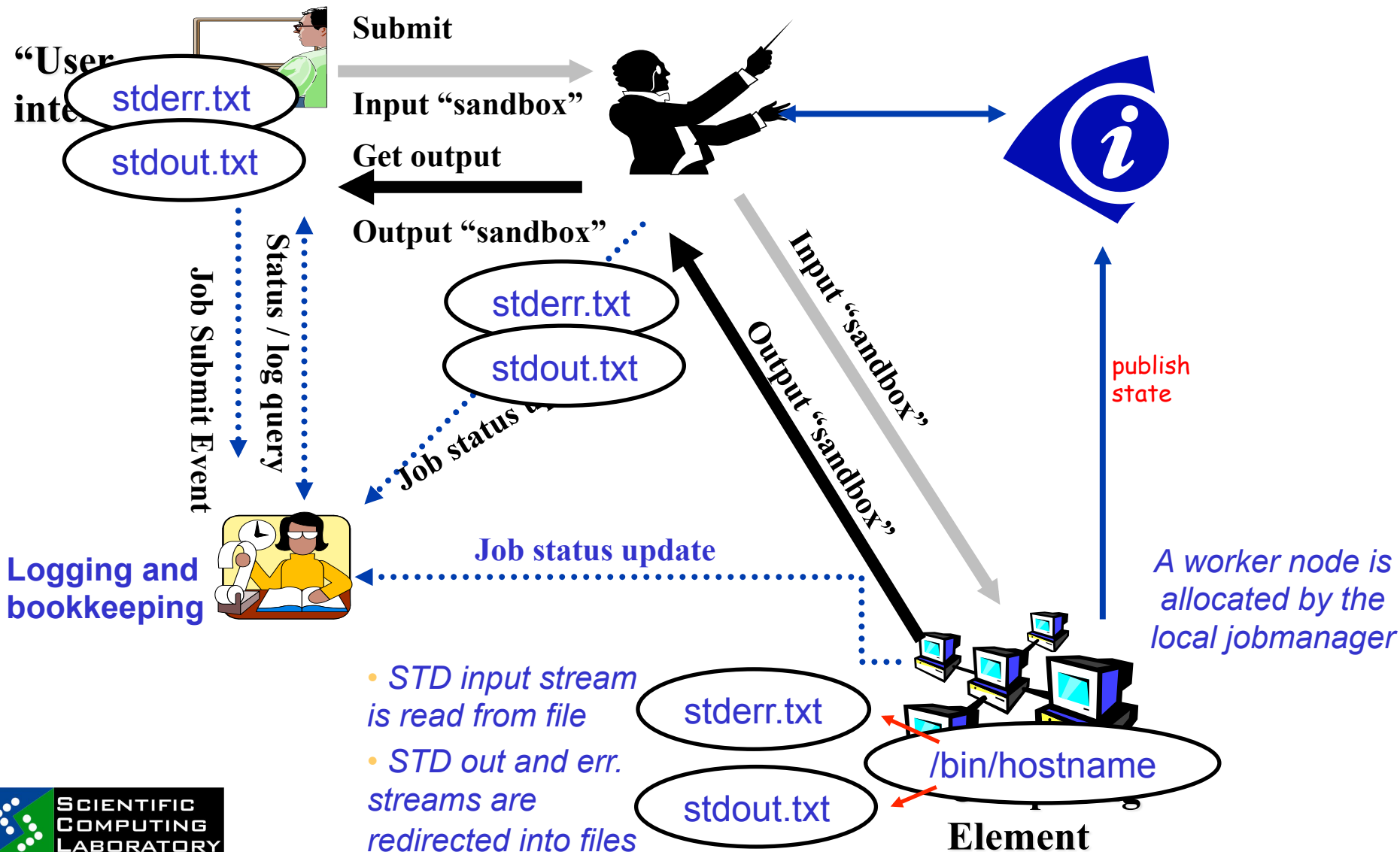
```
Type = "Job";
Executable = "/bin/hostname";
Arguments = "";
StdError = "stderr.txt";
StdOutput = "stdout.txt";
InputSandbox = "";
OutputSandbox = {"stderr.txt", "stdout.txt"};
```

**\$ glite-wms-job-submit -a my.jdl**

**Returns a "job-id" used to monitor job, retrieve output...**



- **Type** – “Job” for sequential jobs; later more details
- **Executable** – sets the name of the executable file;
- **Arguments** – command line arguments of the program;
- **StdOutput, StdError** - files for storing the standard output and error messages output;
- **InputSandbox** – set of input files needed by the program, including the executable;
- **OutputSandbox** – set of output files which will be written during the execution, including standard output and standard error output; these are sent from the CE to the WMS for you to retrieve
- **ShallowRetryCount** – in case of grid error, retry job this many times (“Shallow”: before job is running)



- **Script:**
  - No compilation is necessary
  - Can invoke binary that is statically installed on the CE
- **Binary:**
  - Must be compiled on the User Interface → binary compatibility with CEs is guaranteed
  - Statically linked → to avoid errors caused by library versions
- **Coming from client side**
  - Part of InputSandbox
- **Installed on the CE**
  - Standard software in Linux
  - VO specific software: advertised in information system
- *Use JDL to navigate job to such a site*



WMS version	LCG-2 WMS	gLite WMS via NS gLite 3.0	gLite WMS via WMProxy gLite 3.1+
Delegate proxy		<b>D</b>	<b>glite-wms-job-delegate-proxy</b> -d delegID
Submit	<b>edg-job-submit</b> [-o joblist] jdlfile	<b>glite-job-submit</b> [-o joblist] jdlfile	<b>glite-wms-job-submit</b> [-d delegID] [-a] [-o joblist] jdlfile
Status	<b>edg-job-status</b> [-v verbosity] [-i joblist] jobIDs	<b>glite-job-status</b> [-v verbosity] [-i joblist] jobIDs	<b>glite-wms-job-status</b> [-v verbosity] [-i joblist] jobIDs
Logging	<b>edg-job-get-logging-info</b> [-v verbosity] [-i joblist] jobIDs	<b>glite-job-get-logging-info</b> [-v verbosity] [-i joblist] jobIDs	<b>glite-wms-job-logging-info</b> [-v verbosity] [-i joblist] jobIDs
Output	<b>edg-job-get-output</b> [-dir outdir] [-i joblist] jobIDs	<b>glite-job-get-output</b> [-dir outdir] [-i joblist] jobIDs	<b>glite-wms-job-output</b> [-dir outdir] [-i joblist] jobIDs
Cancel	<b>edg-job-cancel</b> [-i joblist] jobID	<b>glite-job-cancel</b> [-i joblist] jobID	<b>glite-wms-job-cancel</b> [-i joblist] jobID
Compatible resources	<b>edg-job-list-match</b> jdlfile	<b>glite-job-list-match</b> jdlfile	<b>glite-wms-job-list-match</b> [-d delegID] [-a] jdlfile

```

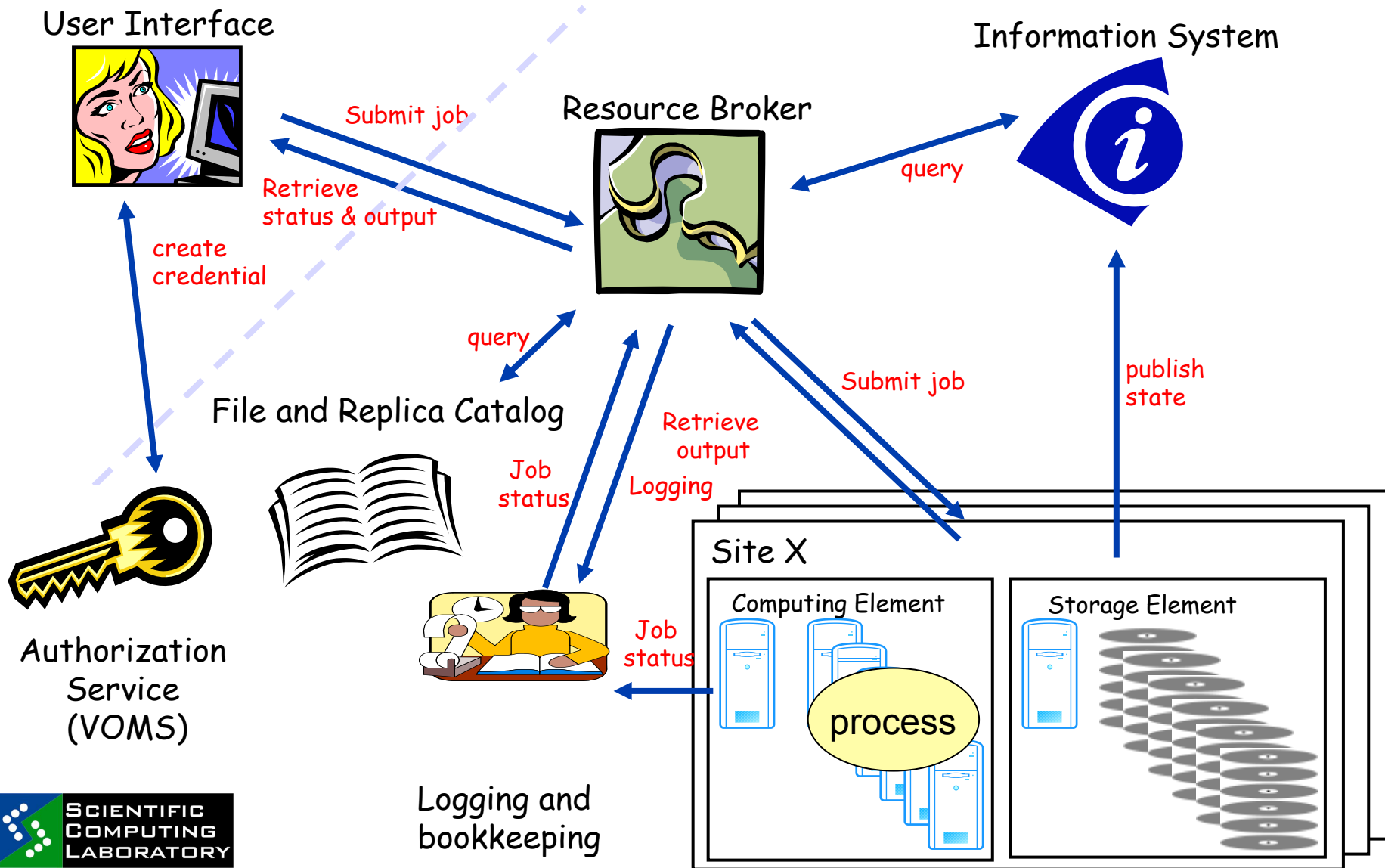
Type = "Job";
Executable = "/bin/hostname";
Arguments = "";
StdError = "stderr.txt";
StdOutput = "stdout.txt";
InputSandbox = "";
OutputSandbox = {"stderr.txt", "stdout.txt"};

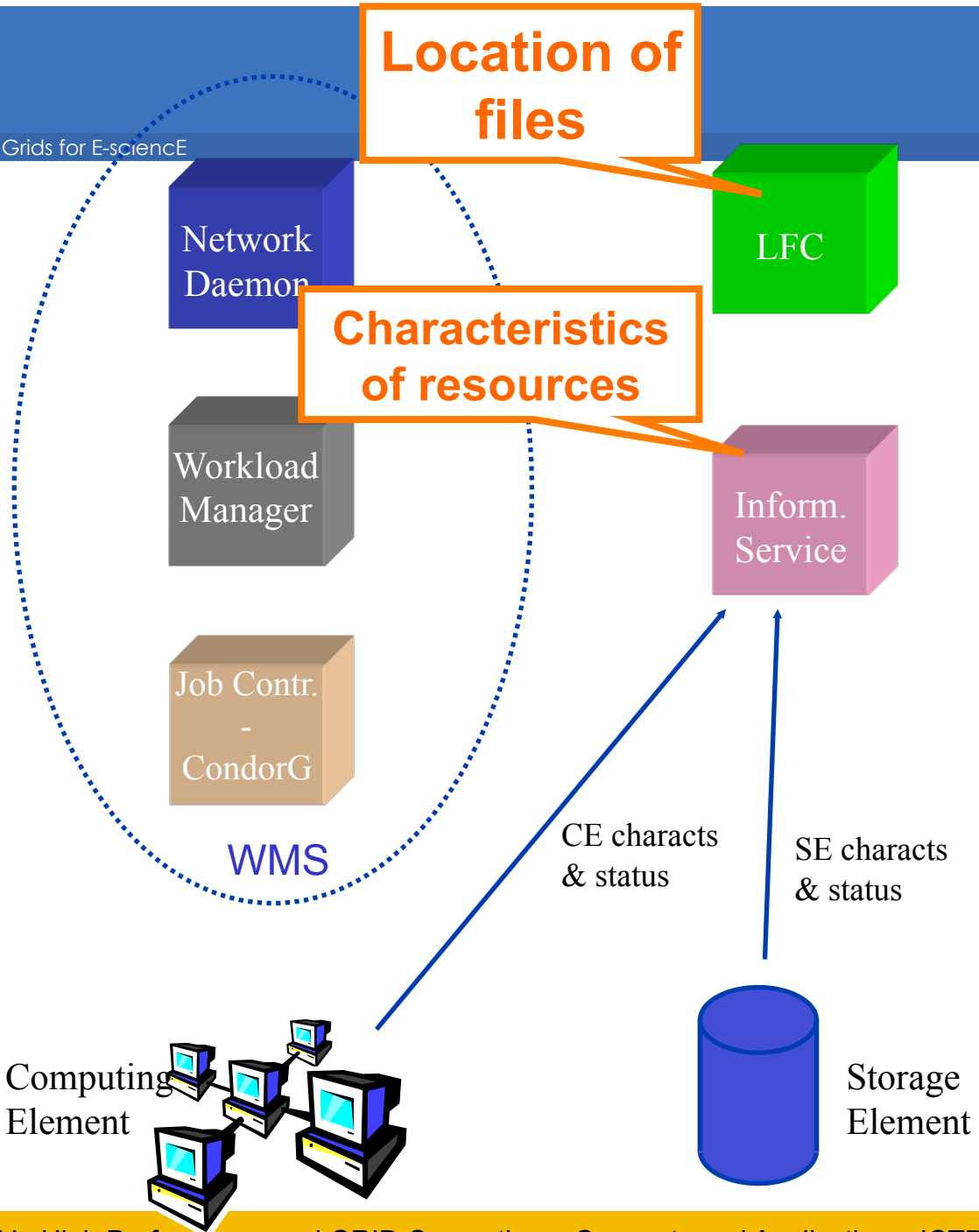
Requirements = other.Architecture=="INTEL" &&
  other.GlueCEInfoTotalCPUs > 480;
Rank = other.GlueCEStateTotalJobs;

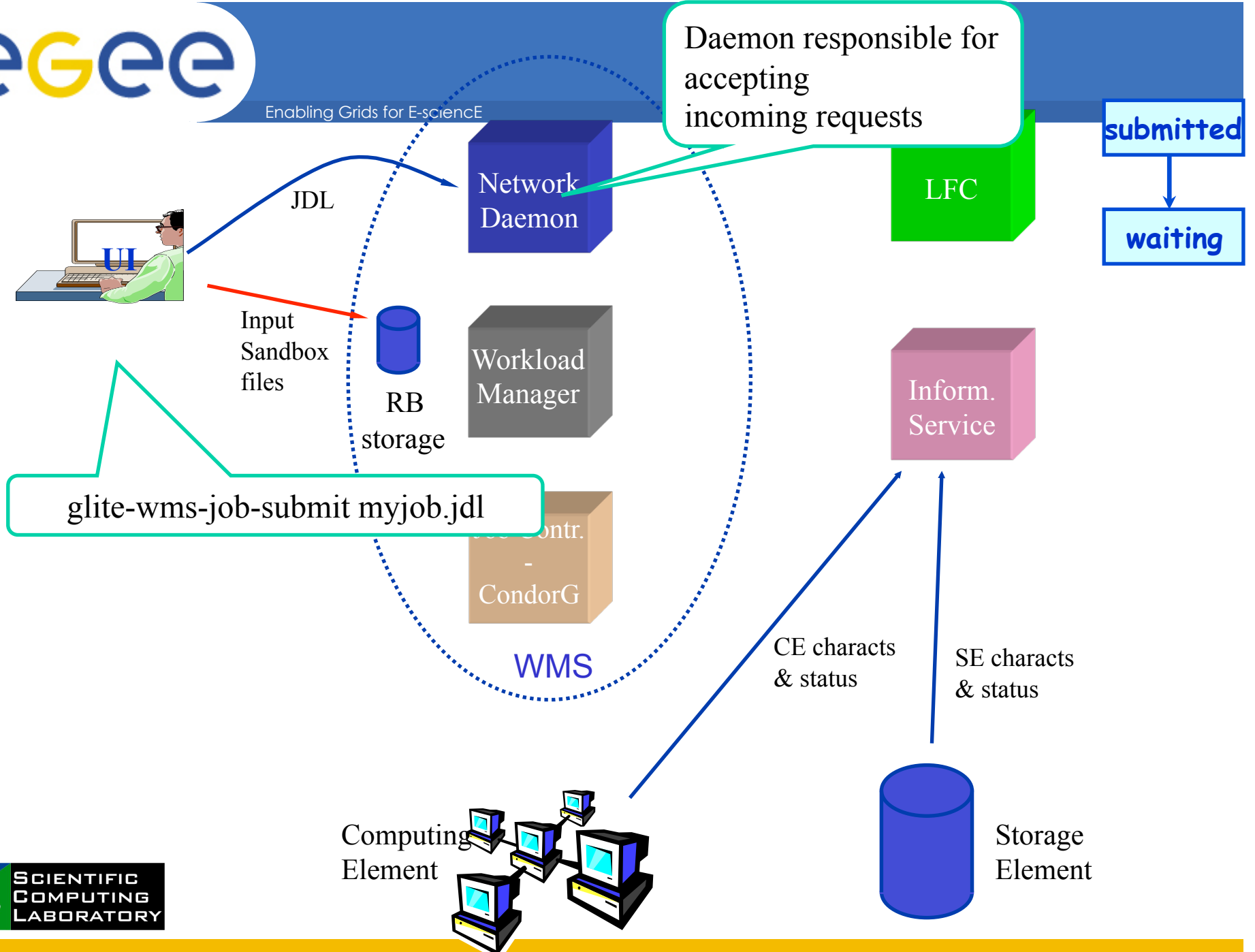
InputData = "lfn:/grid/gridbox/antun/file.txt";
  
```

**WMS uses Information System to find CE**

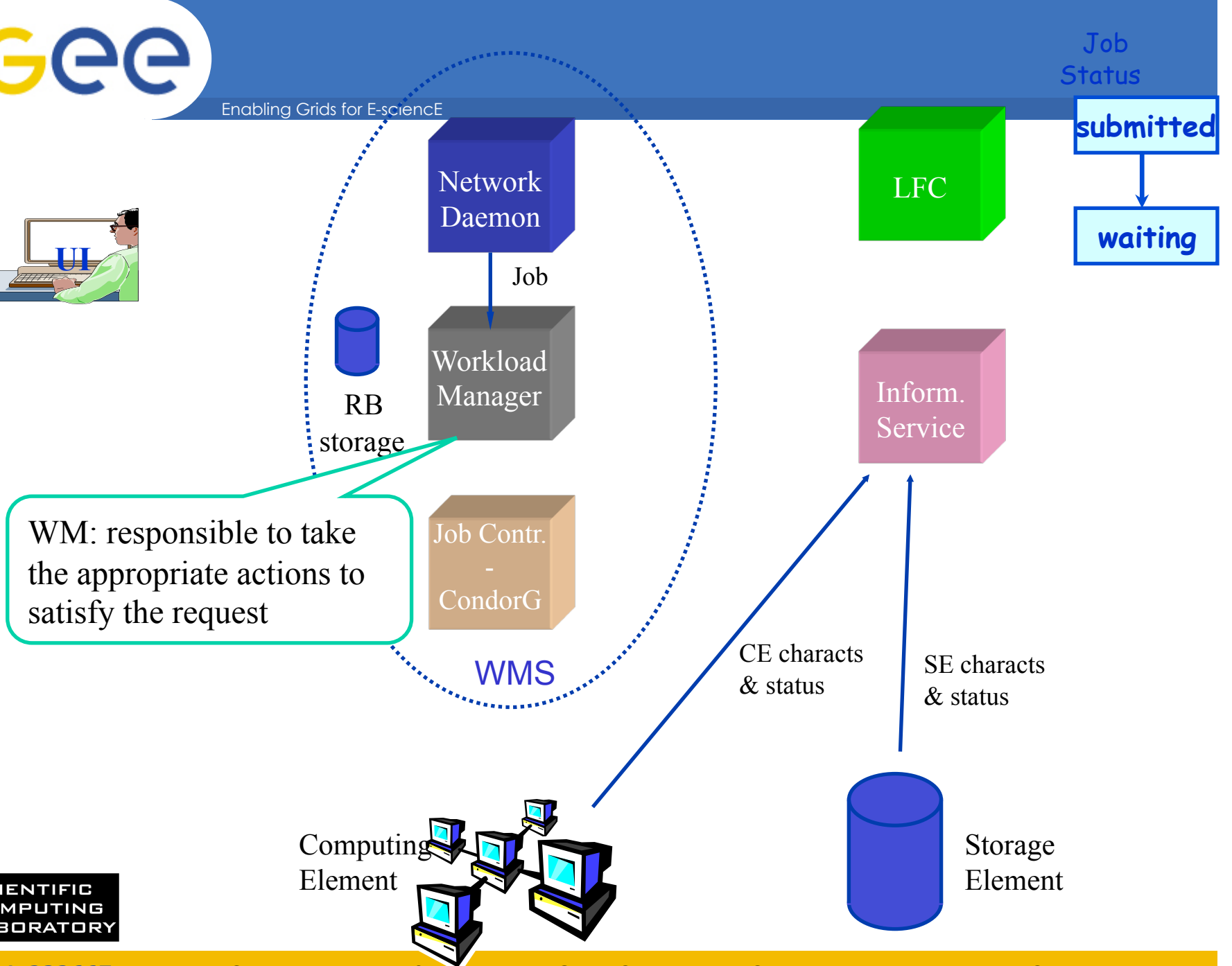
**WMS uses File Catalog to find file location**

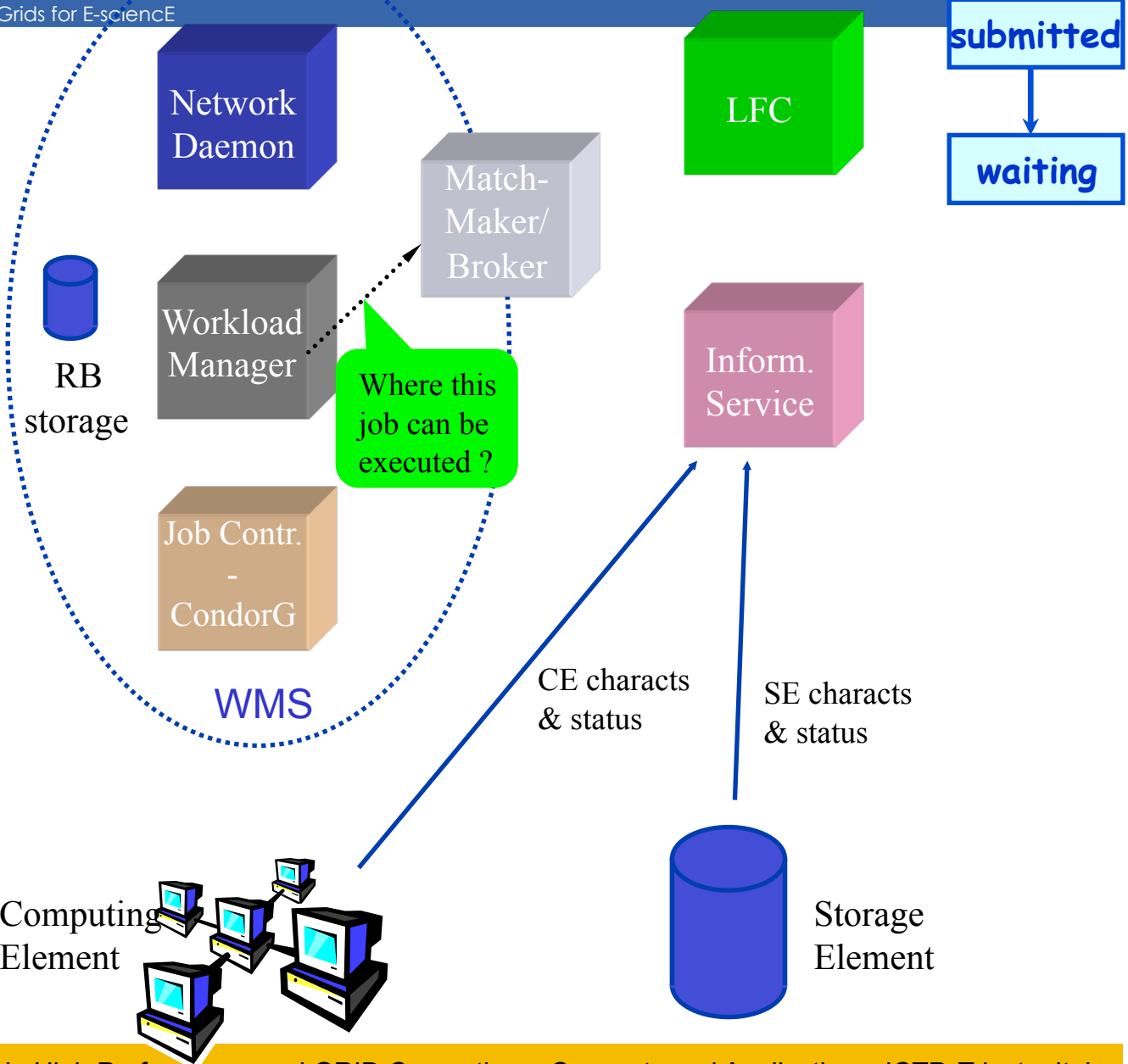






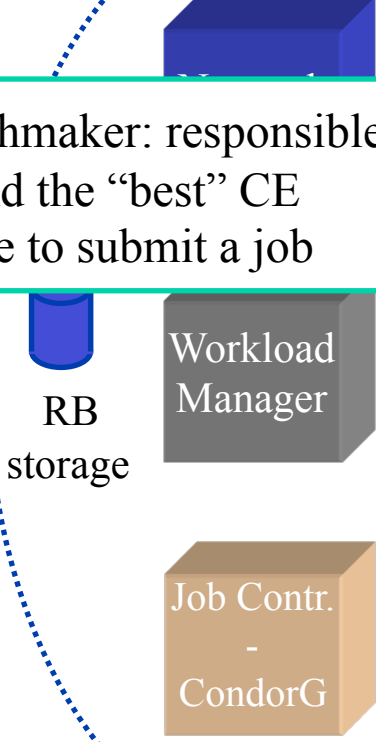








Matchmaker: responsible to find the "best" CE where to submit a job



WMS

Match-Maker/Broker

LFC

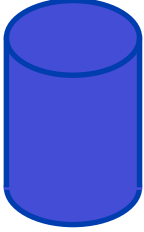
submitted

waiting

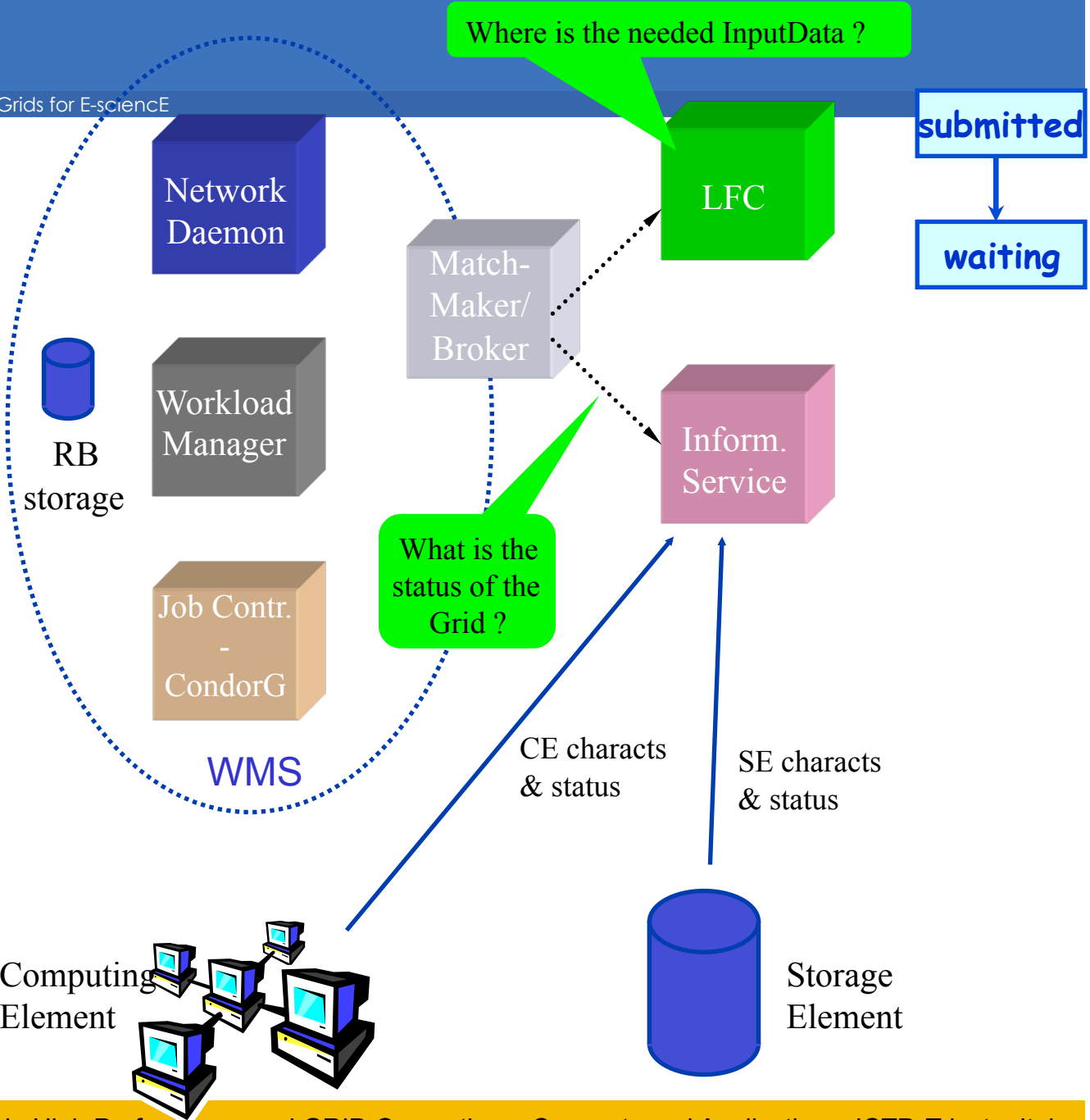
Inform. Service

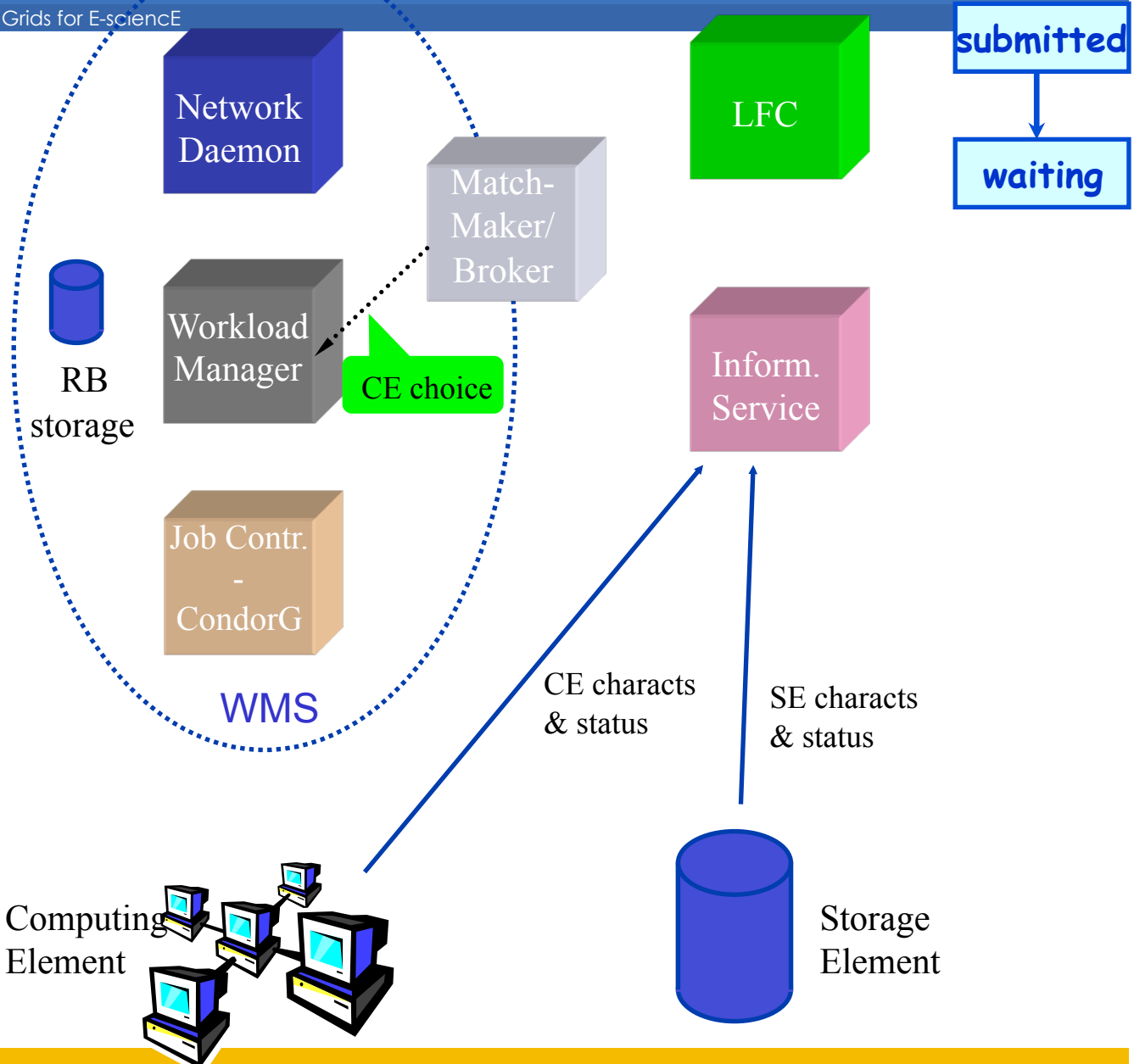
CE characts & status

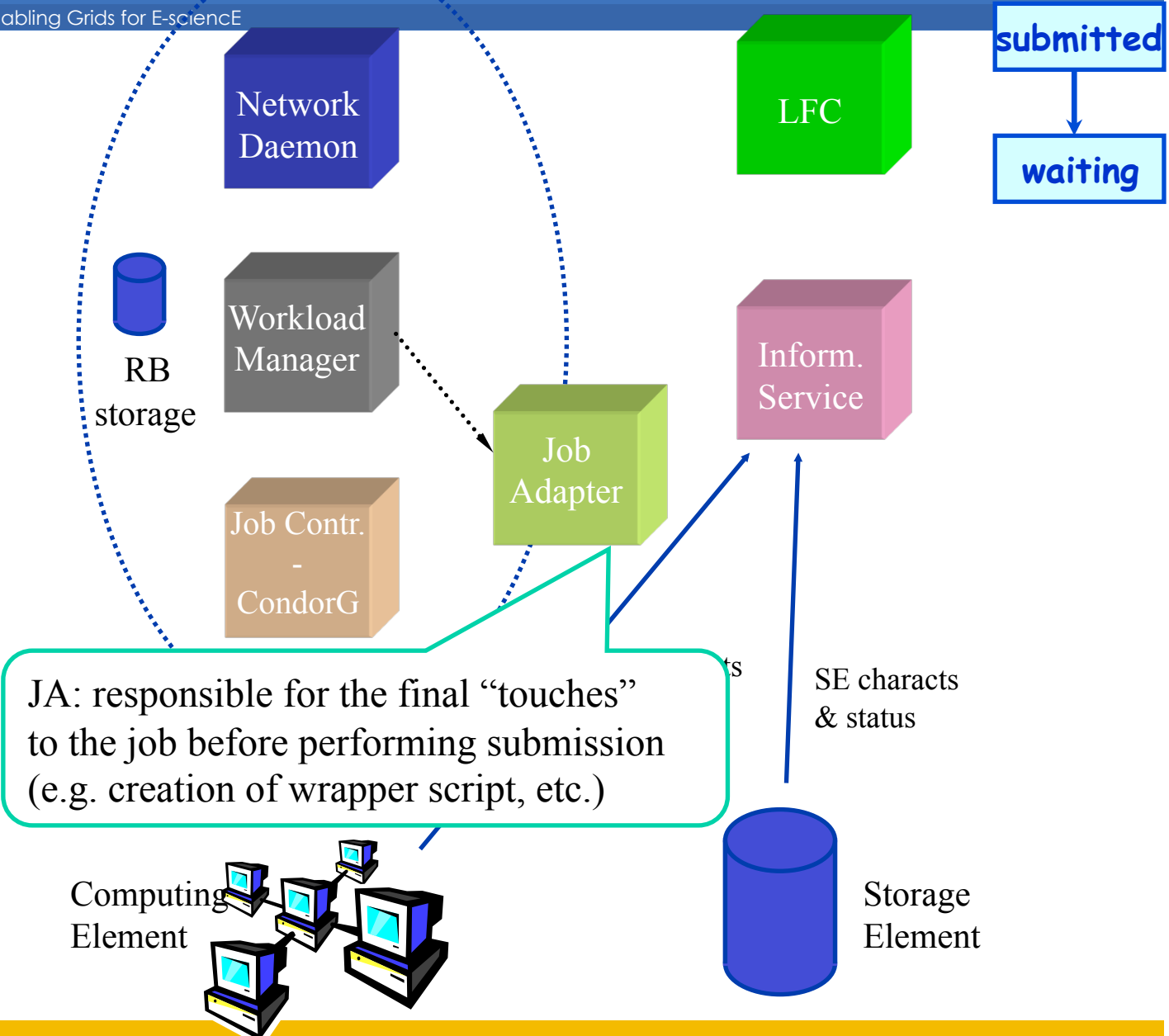
SE characts & status



Storage Element

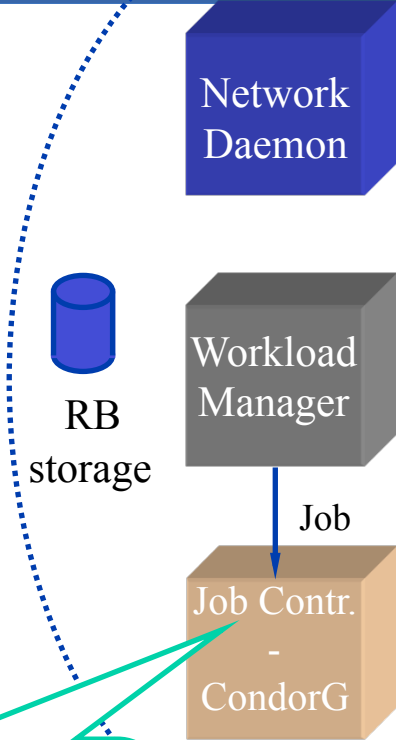








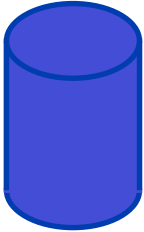
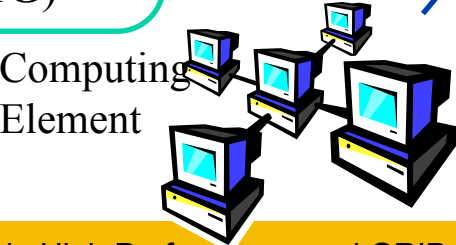
JC: responsible for the actual job management operations (done via CondorG)



submitted

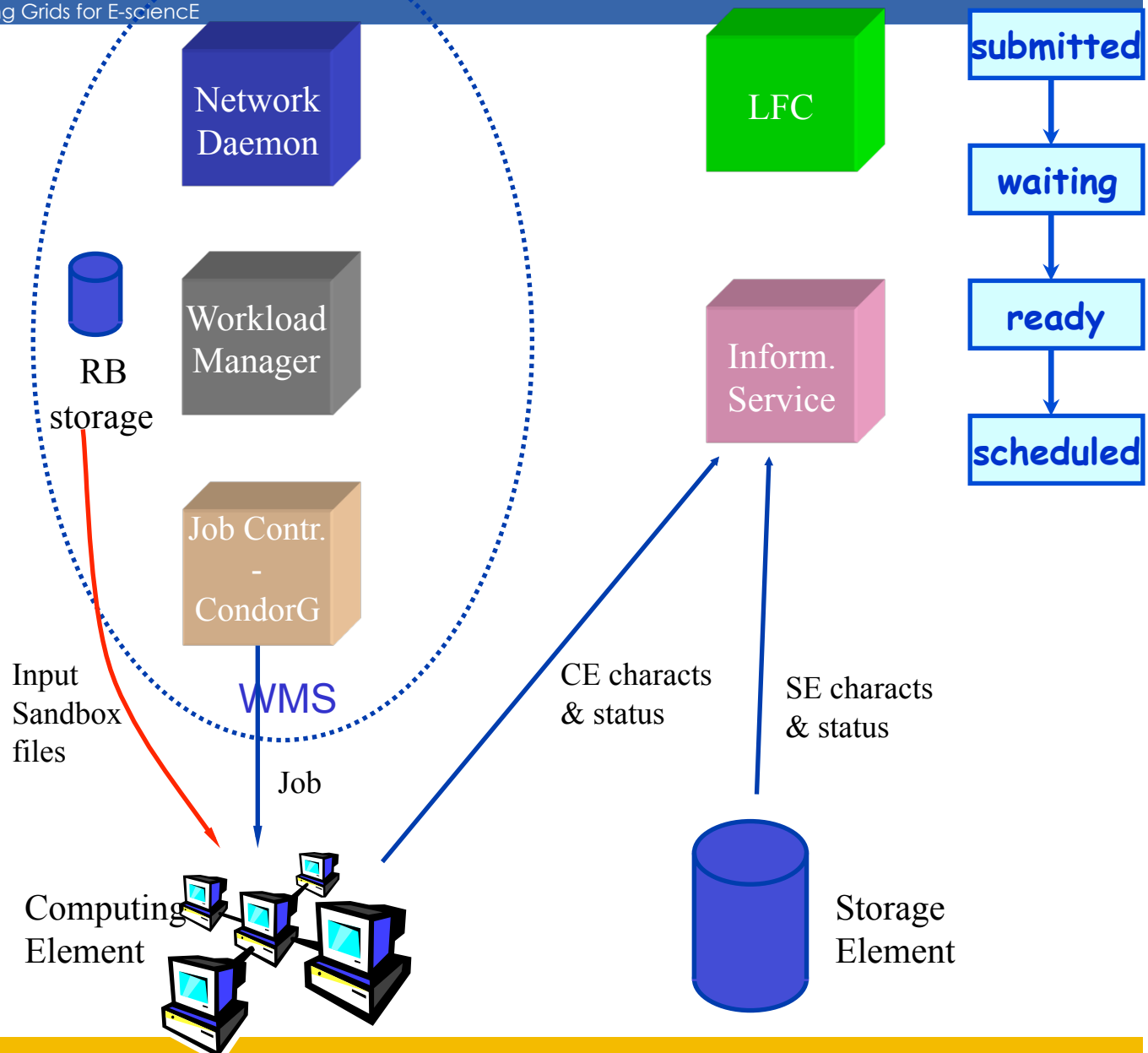
waiting

ready

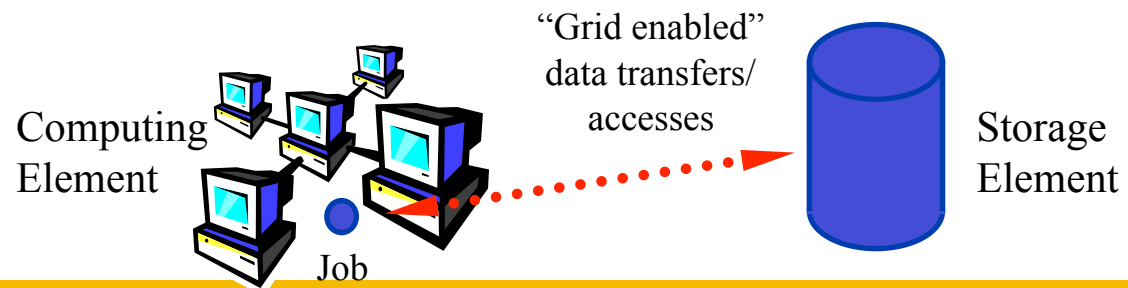
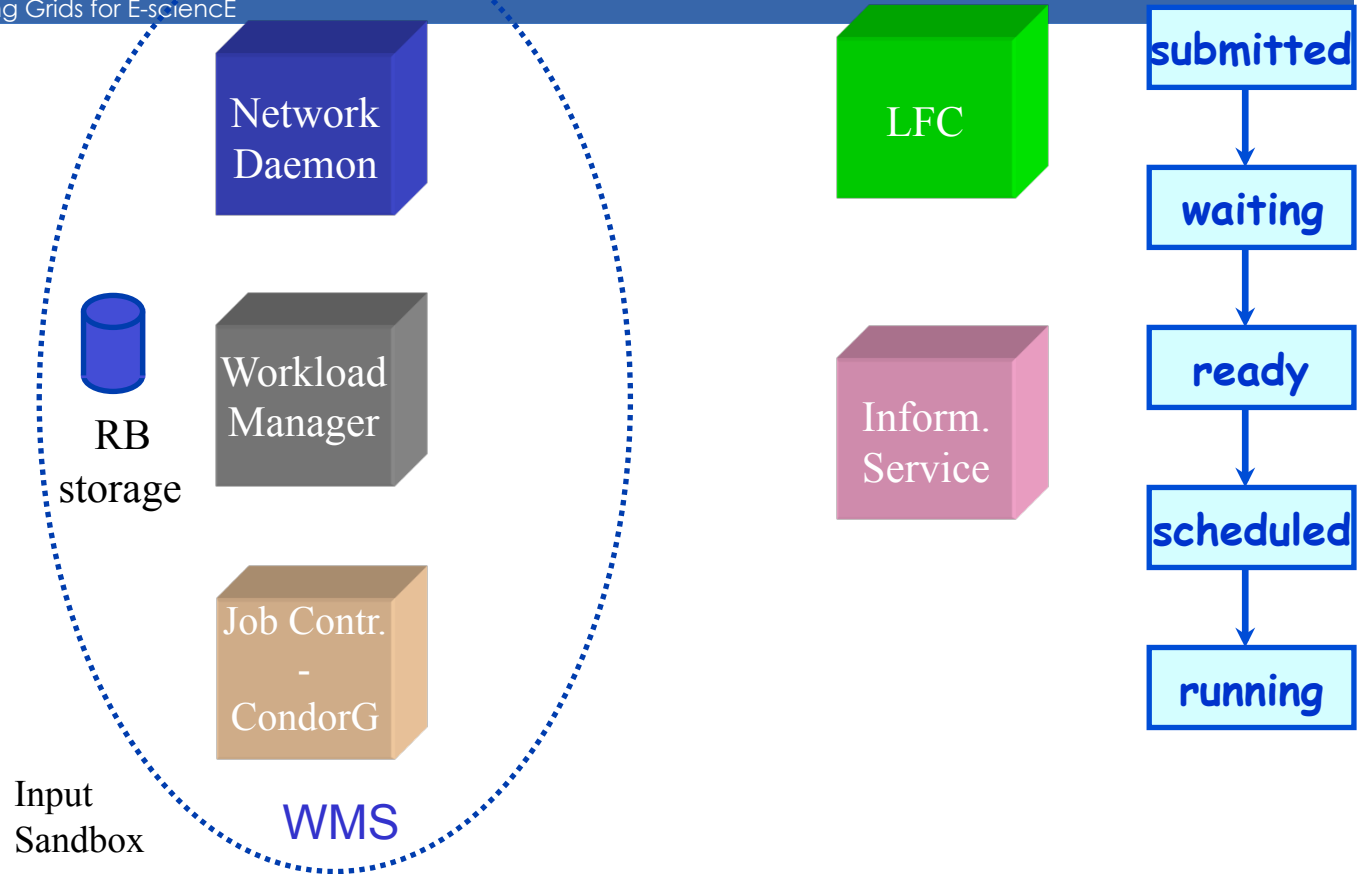


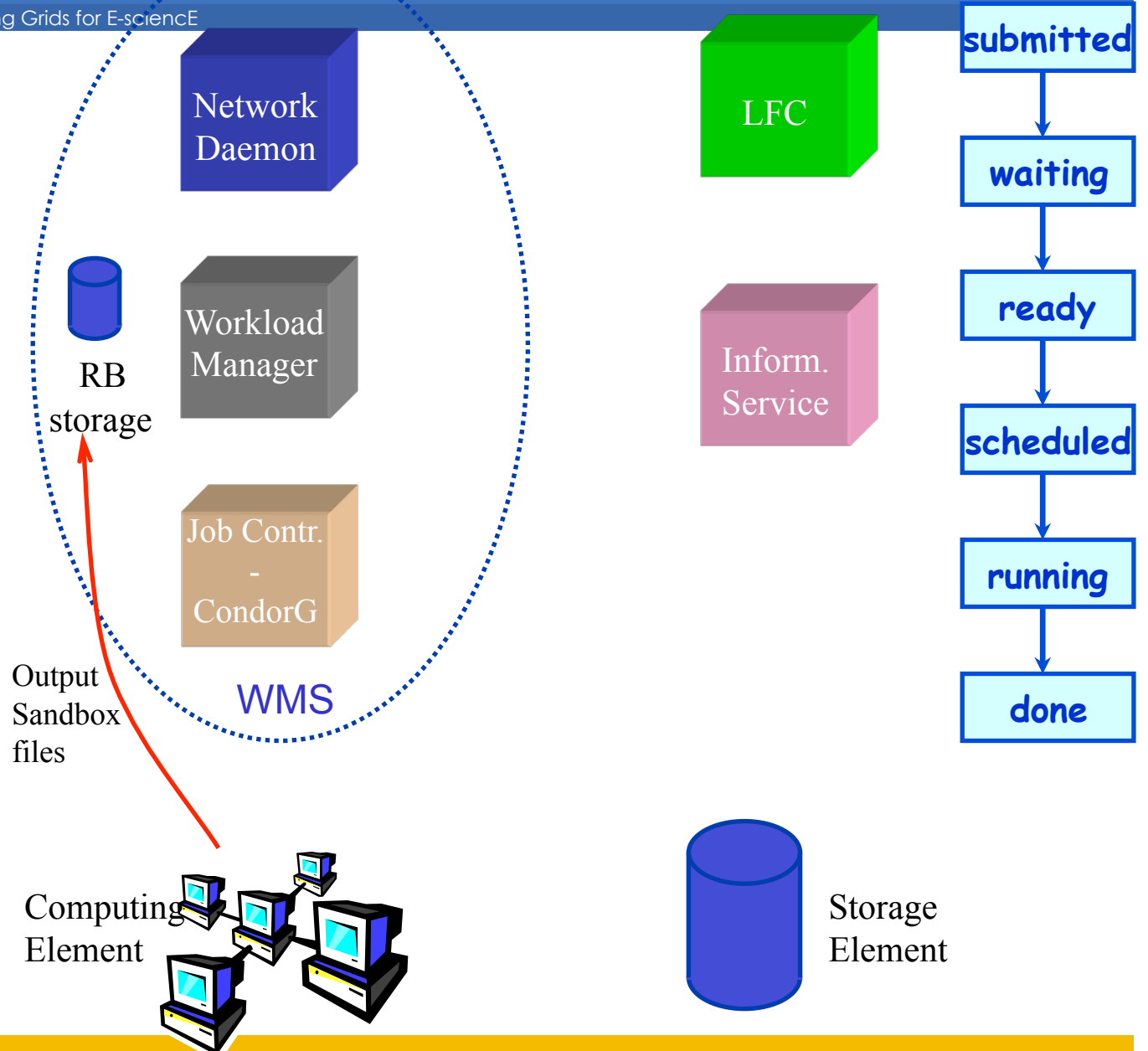
CE characts & status

SE characts & status



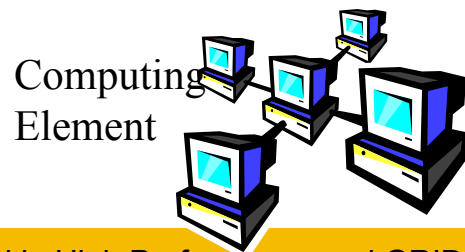
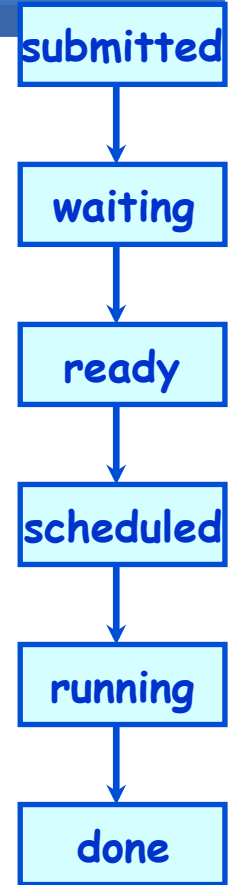
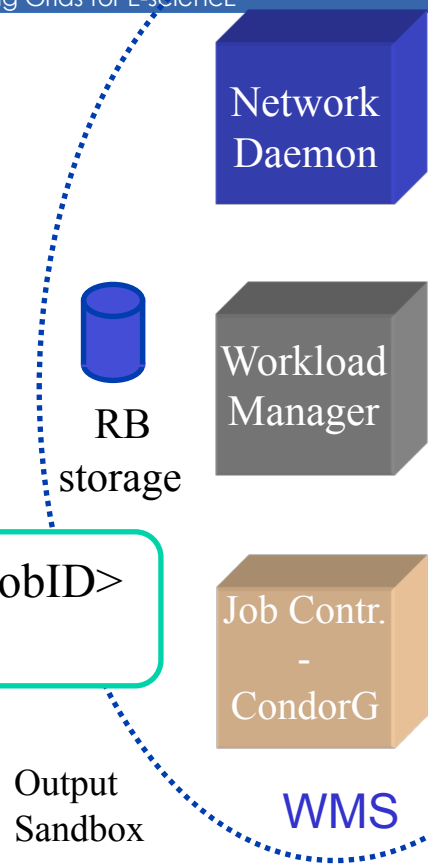








`glite-wms-get-output <jobID>`



Storage Element



Output  
Sandbox  
files



RB  
storage

Network  
Daemon

Workload  
Manager

Job Contr.  
-  
CondorG

WMS

LFC

Inform.  
Service

submitted

waiting

ready

scheduled

running

done

cleared

Computing  
Element



Storage  
Element



# Job monitoring

glite-wms-job-status <jobID>  
glite-wms-job-logging-info <jobID>



LB: receives and stores  
job events; processes  
corresponding job status

Job  
status

Logging &  
Bookkeeping

LB  
proxy

Network  
Daemon

Workload  
Manager

Job Contr.  
-  
CondorG

WMS

Log of  
job events

Computing  
Element



## 1. Meet CE requirements

(defined by Requirements part of JDL)

## 2. Select CE which is close to InputData

- “Close” relationship is defined between CEs and SEs by site administrators
- “Close” is not necessarily physical distance – rather bandwidth
- “Close” typically means same site
  - SE: se-3.grid.box
  - CE: ce-1.grid.box:2119/jobmanager-lcgpbs-gridbox

## 3. Select CE with highest rank

(rank formula is defined by Rank part of JDL)

- **GlueCEUniqueID** – Identifier of a CE  
Eliminating an erroneous CE:  
**other.GlueCEUniqueID !=**  
**"ce-1.grid.box:2119/jobmanager-lcgpbs-gridbox"**
  - **GlueCEInfoTotalCPUs** – max number of CPUs at a CE  
**Rank = other. GlueCEInfoTotalCPUs;**
  - **GlueCEStateWaitingJobs** – number of waiting jobs
  - **GlueCEPolicyMaxCPUTime** – job will be killed after this number of minutes  
**other.GlueCEPolicyMaxCPUTime > 300;**
  - **GlueHostMainMemoryRAMSize** – memory size  
**other.GlueHostMainMemoryRAMSize > 1024;**
- <http://glite.web.cern.ch/glite/documentation/> → JDL specification

- Rank =  
 ( other.GlueCEStateWaitingJobs == 0 ?  
 other.GlueCEStateFreeCPUs : -other.GlueCEStateWaitingJobs);

## if there are no waiting jobs,

- then the selected CE will be the one with the most free CPUs
- else the one with the least waiting jobs.

- Requirements =  
 ( Member(„IDL2.1”,  
 other.GlueHostApplicationSoftwareRunTimeEnvironment) ) &&  
 (other.GlueCEPolicyMaxWallClockTime > 10000);

## CE where,

- IDL2.1 software is available
- At least 10000s can be spent on the site (waiting + running)

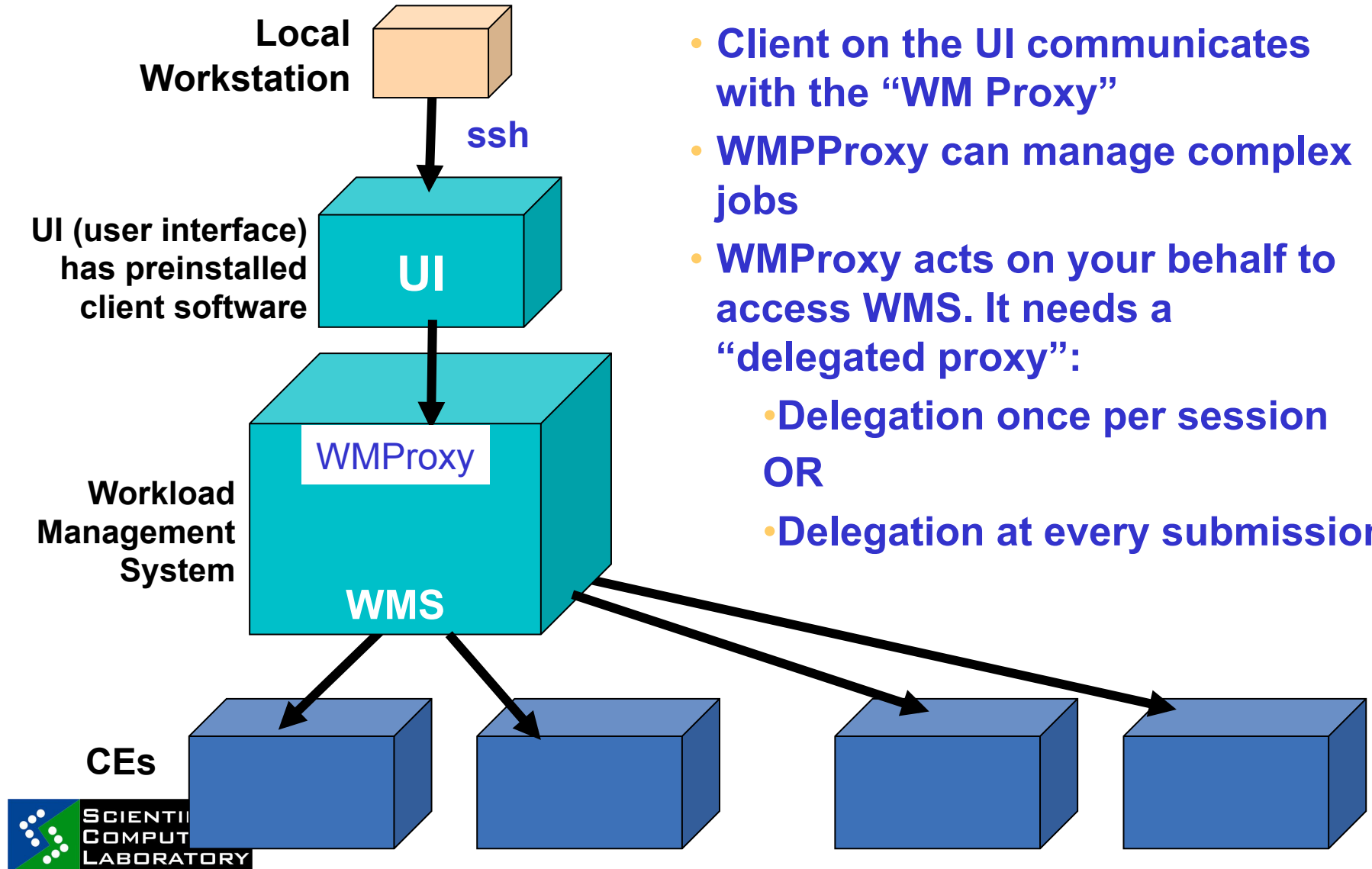


- `other.GlueHostMainMemoryRAMSize >=`  
`512 * ( other.GlueHostArchitectureSMPSize > 0 ?`  
`other.GlueHostArchitectureSMPSize : 1 )`

**At least 512 MB of RAM memory per CPU core should be available**

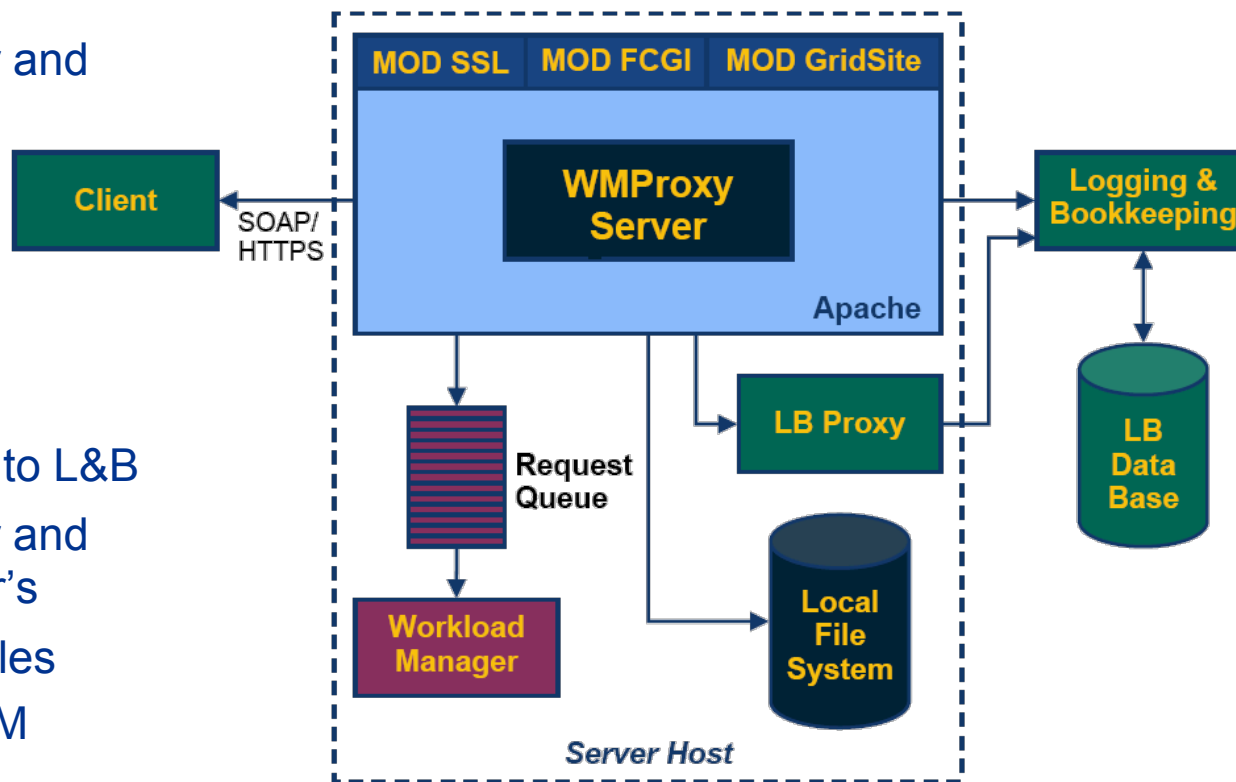
- `other.GlueHostArchitecturePlatformType == "x86_64"`  
**x86\_64 arch requested**

Flag	Meaning
SUBMITTED	submission logged in the Logging & Bookkeeping service
WAIT	job match making for resources
READY	job being sent to executing CE
SCHEDULED	job scheduled in the CE queue manager
RUNNING	job executing on a Worker Node of the selected CE queue
DONE	job terminated without grid errors
CLEARED	job output retrieved
ABORT	job aborted by middleware, check <i>reason</i>



- Client on the UI communicates with the “WM Proxy”
- WMPProxy can manage complex jobs
- WMPProxy acts on your behalf to access WMS. It needs a “delegated proxy”:
  - Delegation once per session
  - OR
  - Delegation at every submission

- **WMPProxy is a SOAP Web service providing access to the Workload Management System (WMS)**
- **Job characteristics specified via JDL**
  - jobRegister
    - create id
    - map to local user and create job dir
    - register to L&B
    - return id to user
  - input files transfer
  - jobStart
    - register sub-jobs to L&B
    - map to local user and create sub-job dir's
    - unpack sub-job files
    - deliver jobs to WM



- **Type**
  - Job, DAG, Collection
- **JobType (only when the Type is set to Job!)**
  - Normal (sequential batch job), Parametric, Interactive, MPICH, ~~Checkpointable, Partitionable~~
- **Executable**
  - The name of the executable (absolute path)
- **Arguments**
  - Job command line arguments
- **StdInput, StdOutput, StdError**
  - Standard input/output/error of the job (stdin absolute path; stdout & stderr relative path)
- **Environment**
  - List of environment variables to be set for the binary
- **InputSandbox**
  - List of files on the UI local disk needed by the job for running
  - The listed files will be staged to the remote resource
- **OutputSandbox**
  - List of files, generated by the job, which have to be retrieved
  - Files will be transferred back

- **Input Data**
  - For the broker, WMS does not transfer these files
- **Output Data**
  - For the broker, WMS does not transfer these files
- **Requirements**
  - Required CE characteristics
- **Rank**
  - “Goodness” value for compatible CEs
- **ShallowRetryCount**
  - in case of grid error, retry job this many times
  - “Shallow”: before job is running
- **RetryCount**
  - resubmit if the job failed in Running mode
  - If job fails after it has already done something (e.g. creating a Grid file) then resubmission can generate inconsistencies
- **MyProxyServer**
  - where to download proxy from in case of the existing proxy expires
  - Done by WMS

- **A set of independent jobs**
- **For some reason must be managed as a single unit**
- **Possible reasons:**
  - Belong to the same experiment
  - Share common input files
  - Optimize network traffic
- **Sharing of sandboxes**

[

**Type = "Collection";**

*Transfer from UI only once*

```
InputSandbox = {
    "sharedFile1"; . . . ; "sharedFileM" };
```

**nodes = {**

```
[ JobType = "Normal";
  InputSandbox = {root.InputSandbox, . . . }
  ...; ],
```

*JDL of 1<sup>st</sup> job*

. . .

```
[ JobType = "Normal";
  ...; ],
```

*JDL of N<sup>th</sup> job*

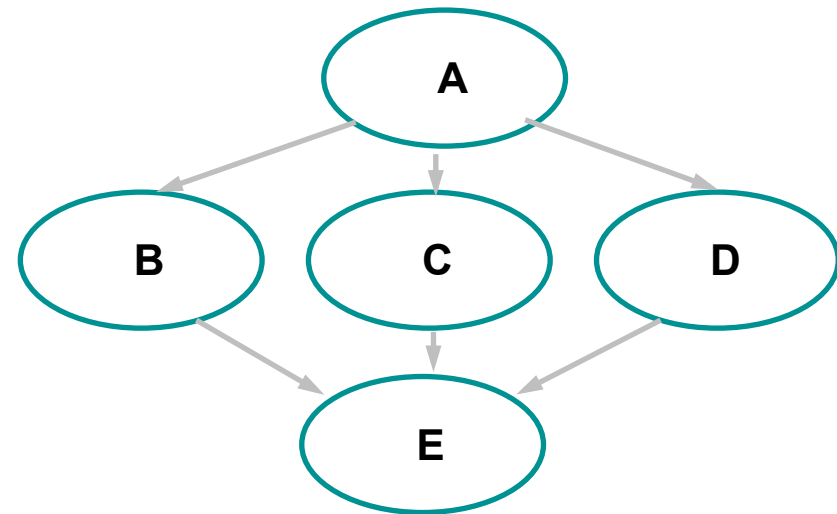
. . .

```
};
```

]



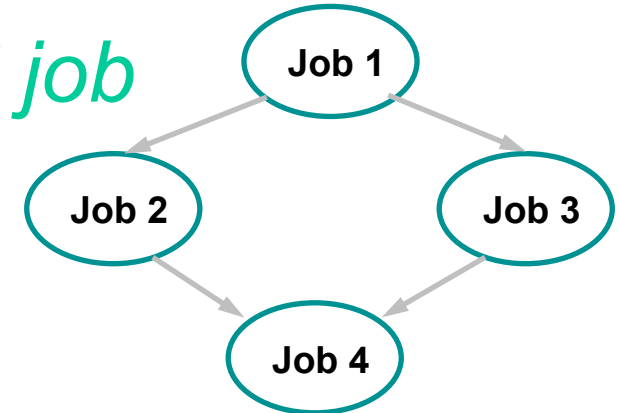
- **Direct Acyclic Graph (DAG)** is a set of jobs where the input, output, or execution of one or more jobs depends on one or more other jobs
- Sharing and inheritance of sandboxes
  - **Include OutputSandbox in the next InputSandbox**
- Dependencies defined between pairs of jobs



```
[ Type = "DAG";
  InputSandbox = {
    "sharedFile1", . . . , "sharedFileN" };
  nodes = [
    job1 = [
      description = [
        JobType = "Normal";
        . . . ; ] ;
    . . .
  ] ;
  dependencies = {
    {job1, {job2, job3}}, {job2, job4}, {job3, job4} };
  } ;
]
```

*Transfer from UI only once*

*JDL of 1<sup>st</sup> job*



*Graph structure*

```
[ Type = "DAG";
...
  job4 = [
    description = [
      JobType = "Normal";
      InputSandbox = {
root.nodes.job1.description.OutputSandbox[0],
root.nodes.job2.description.OutputSandbox,
...};
...];
];
]
```

- A set of jobs generated from one JDL
- Useful where many similar (but not identical) jobs must be executed
  - Parameter study, parametric sweep applications
  - Majority of grid applications are parametric!
- One or more parametric attributes in the JDL:
  - Use the `_PARAM_` keyword
  - E.g. `InputSandbox = "input_PARAM_";`

```
[  
  Type = "Parametric";  
  . . .  
  
  ParameterStart = 0;  
  ParameterStep = 2;  
  Parameters = 6;    → _PARAM_: 0, 2, 4, 6, 8, 10  
  
  Arguments = "inputfigure_PARAM.jpg";  
  StdOutput = "transformed_PARAM.jpg";  
  OutputSandbox = {" transformed_PARAM.jpg ",...};  
  . . .  
]
```

```
[  
  Type = "Parametric";  
  . . .  
  
  Parameters = {alpha, beta, gama};  
  
  Arguments = "inputfigure_PARAM_.jpg";  
  StdOutput = "transformed_PARAM_.jpg";  
  OutputSandbox = {" transformed_PARAM_.jpg ",...};  
  . . .  
]
```

- For simple jobs: **glite-wms-...** is the recommended way to use the WMS
- History:
  - Before the **glite-wms-** commands we had **glite-** commands
    - used the WMS without WMPProxy
  - Before the **glite-** commands we had
    - **edg-** commands (edg-job-submit....)
      - *European Data Grid – project before EGEE*
    - Used the “resource broker”
    - Still very widely used
  - You might see these commands still in use.
- Status
  - Complex jobs with WMPProxy: first stable version just released. Not yet in routine production use
  - Watch for news!

- Create and submit a JDL file with different requirements and rankings
- Create and submit a JDL file for a collection of jobs
- Create and submit a JDL file for a parametric job
- Create and submit a JDL file for a DAG job