

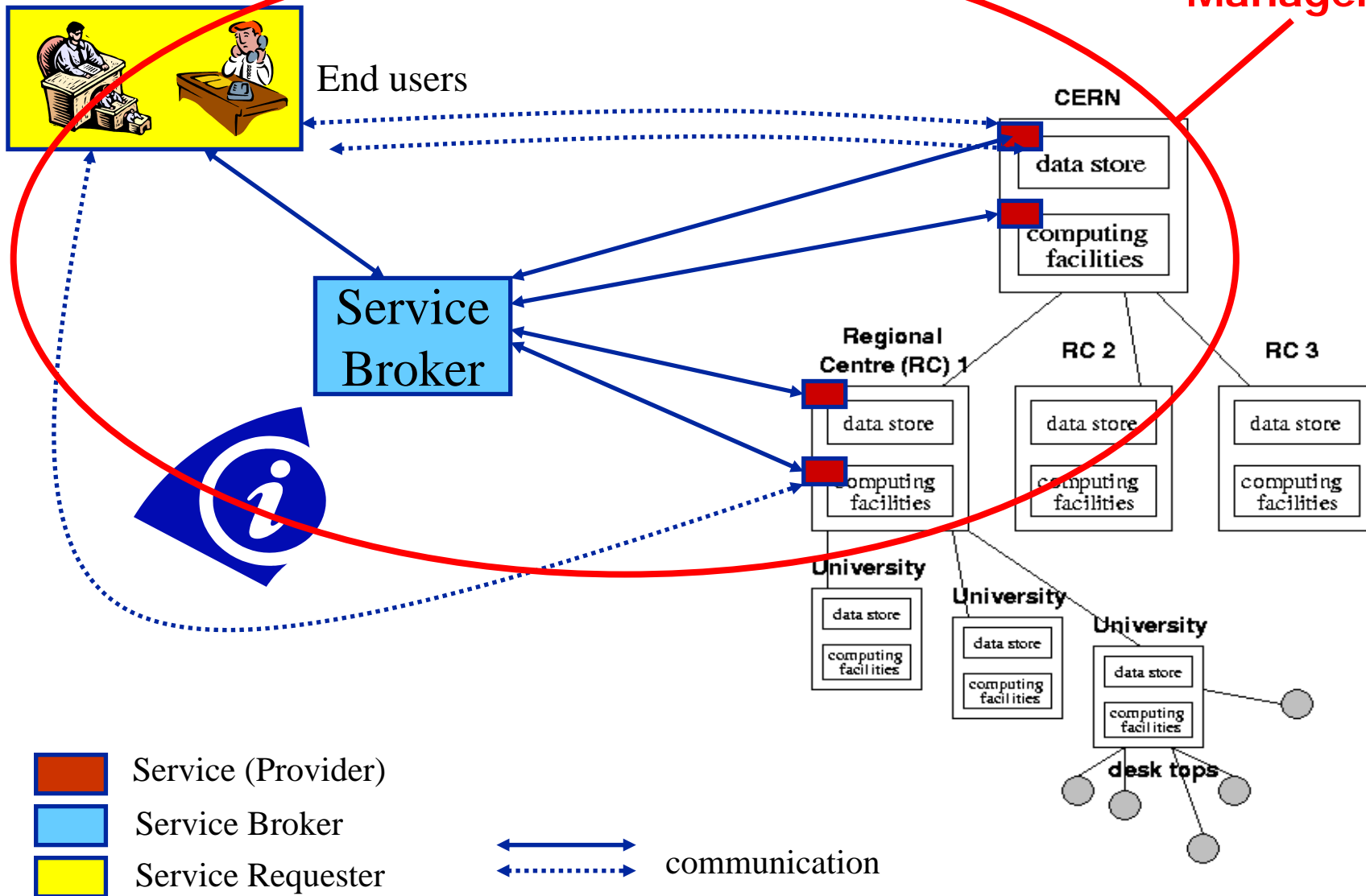
Workload Management



Heinz Stockinger
CERN & INFN

A Reference Grid

Workload Management



Contents

General concepts of Grid Workload Management Systems

EGEE-0 Workload Management Systems

Job Preparation

Architecture / Job submission and status monitoring

Matchmaking

Different job types

Workload Management - Definitions

- ◆ Resource Management includes the **efficient usage** of **computing and storage** resources
 - Processor time, memory, storage, network, etc.
- ◆ Here, mainly referred to as “**workload management system**” since it deals with the distribution of user executables to Grid resources
 - Don't put the load on one subsystem but distribute it
 - Manage the workload produced by end users
 - Workload management is partly referred to as “**scheduling**”
- ◆ From the user's point of view, workload management should be transparent
- ◆ Workload consists of **user jobs**
- ◆ A **job** can be **any kind of executable that requires CPU or storage**

General concepts of Grid Workload Management Systems

EGEE-0 Workload Management Systems

Job Preparation

Architecture / Job submission and status monitoring

Matchmaking

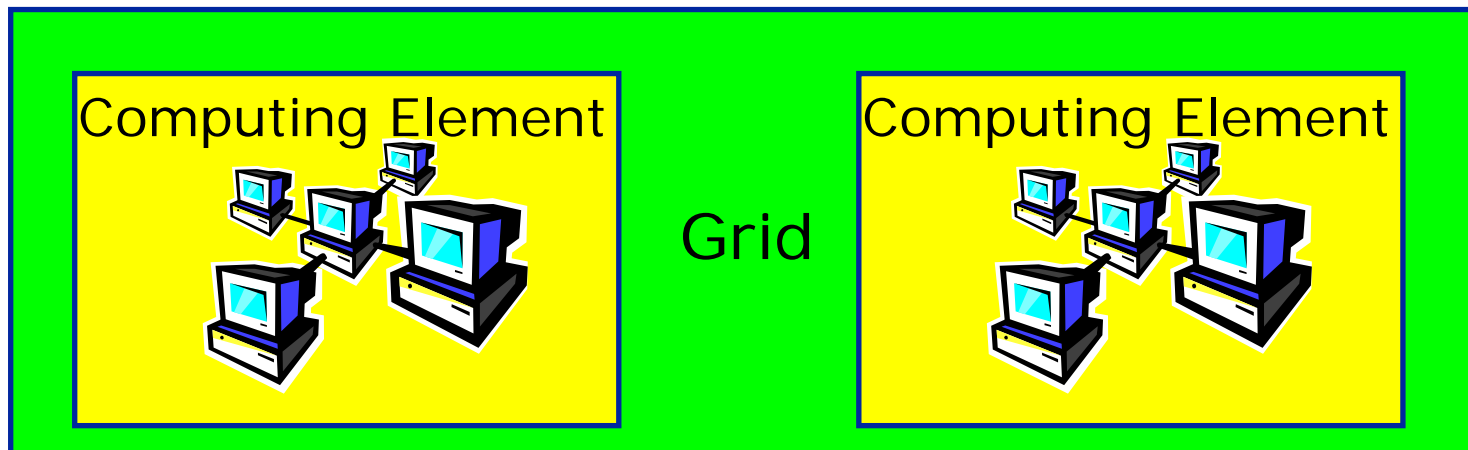
Different job types

Grid Workload Management System

- ◆ Scheduling consists of:
 - Resource **Discovery/Brokering**
 - Find suitable resources
 - **Matchmaking**
 - Assign a job to a resource that satisfies job requirements
 - **Job execution**
 - Execute the jobs and retrieve output
 - Deal with error management
- ◆ Needs to interact with all major Grid services/components
 - Information System,
 - Grid Security
 - Computing Element
 - Storage Element
- ◆ Job execution requires to find the “right” Computing Element (computing resource)
 - Main focus here to **execute CPU intensive jobs**

How to allow for execution of CPU intensive jobs?

- ◆ Use parallel and distributed computing environments
 - Parallel machines
 - Clusters (set of workstations)
 - Sometimes also referred to as “farms” if they are loosely coupled
- ◆ Grid = set of clusters or parallel machines
- ◆ Main focus here on **clusters/farms** and **high throughput computing**
 - Rather than optimising the execution of a single program, allow for several user jobs that then efficiently use existing resources

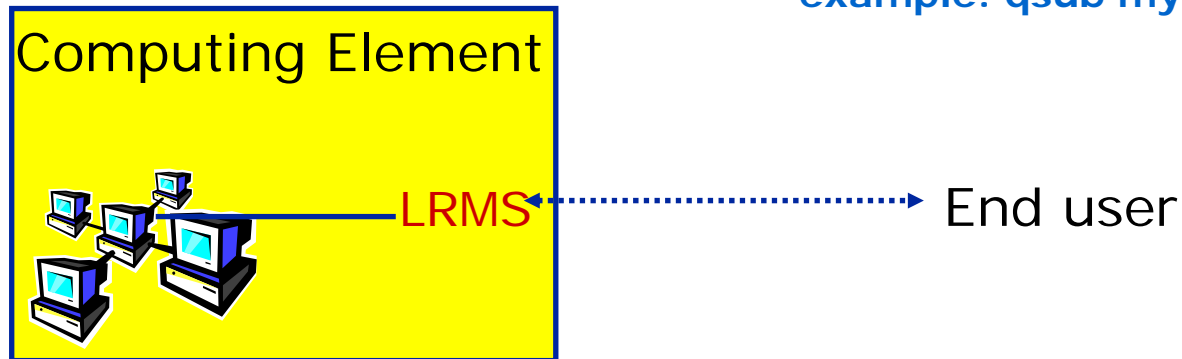


Local Resource Management System (LRMS)

- ◆ Manage the local computing resources in the Computing Element
- ◆ Often, Batch Systems are used
 - PBS (Portable Batch System)
 - LSF (Load Sharing Facility)

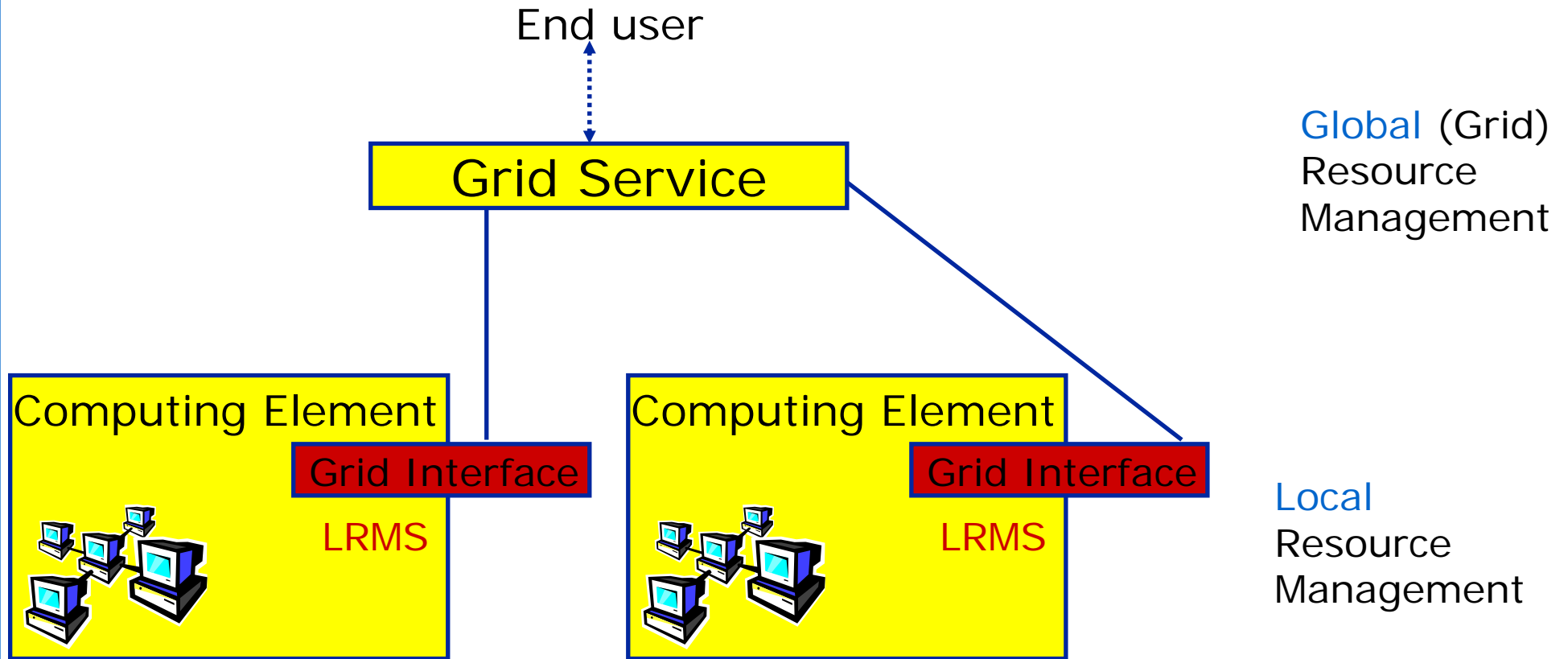
Interaction with local batch system
through native batch client

example: `qsub my.executable`



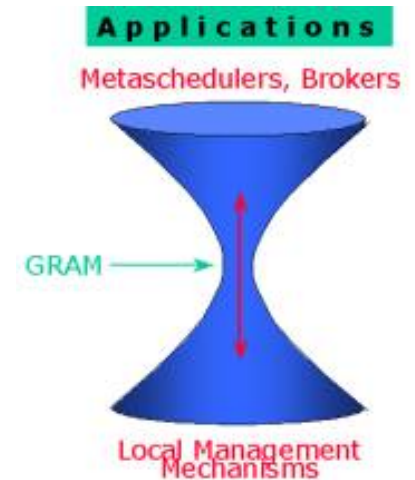
Grid Resource Management

- ◆ Manage several Computing Elements



Globus Resource Management

- ◆ GRAM (Grid Resource Allocation Manager)
- ◆ Service that provides a Grid Interface to Local Resource Management System
 - Also referred to as "Gatekeeper"
 - Provides a general interface to different batch systems like PBS, LSF
 - User needs to specify the exact **hostname of Computing Element**



End user

example: `globus-job-submit host1.cern.ch /bin/lis`

Computing Element

Grid Interface

LRMS



Computing Element

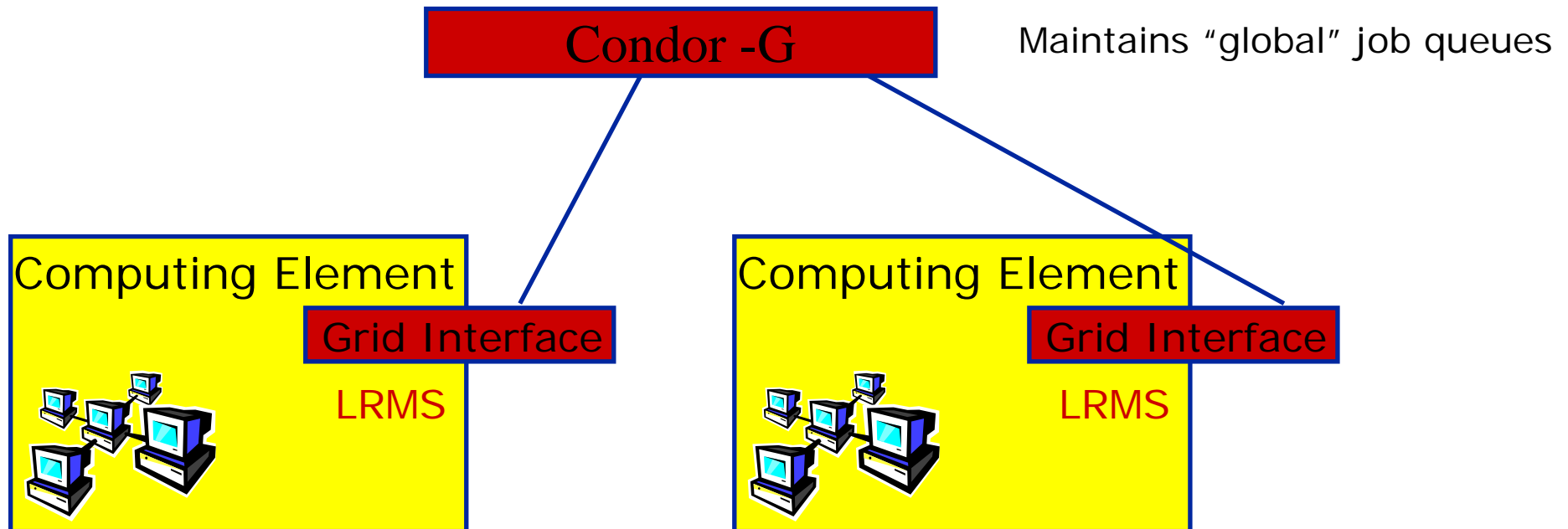
Grid Interface

LRMS



Condor-G

- ◆ From Condor High-throughput computing project
 - <http://www.cs.wisc.edu>
- ◆ Global Resource Management allows to submit to Globus GRAM managed resources
- ◆ No match-making/brokering



The EGEE-0 Workload Management System

- ◆ Builds on both (Globus and Condor) and provides a full Grid workload management system
- ◆ The Goal of WMS is the **distributed scheduling and resource management in a Grid environment.**
- ◆ What does it allow Grid users to do?
 - To submit their jobs
 - To execute them on the “best resources”
 - The WMS tries to optimize the usage of resources
 - To get information about their status
 - To retrieve their output

General Scheduling Approaches

- ◆ Also discussed in the Global Grid Forum area "Scheduling and Resource Management"
 - Generalise architecture, super-scheduling, protocols, Grid economic, job description, ...
- ◆ Scheduling can be categorised as follows:
 - Scheduler organisation
 - ◆ Centralised, distributed, hierarchical
 - Scheduling policy
 - ◆ System-oriented (performance) – application-oriented (throughput)
 - State estimation technique
 - ◆ Non-predictive - predictive

General concepts of Grid Workload Management Systems

EGEE-0 Workload Management Systems

Job Preparation

Architecture /Job submission and status monitoring

Matchmaking

Different job types

Job preparation

- ◆ **Information** to be specified **when a job has to be submitted**:
 - Job characteristics
 - Job requirements and preferences on the computing resources
 - Also including software dependencies
 - Job data requirements
- ◆ Information specified using a Job Description Language (JDL)
 - Based upon *Condor's CLASSified ADvertisement language (ClassAd)*
 - Fully extensible language
 - A ClassAd
 - Constructed with the classad construction operator []
 - It is a sequence of attributes separated by semi-colons.
 - An attribute is a pair (key, value), where value can be a Boolean, an Integer, a list of strings, ...
 <attribute> = <value>;
- ◆ So, the JDL allows definition of a set of attribute, the WMS takes into account when making its scheduling decision

Job Description Language (JDL)

- ◆ The supported attributes are grouped in two categories:
 - **Job Attributes**
 - Define the job itself
 - **Resources**
 - Taken into account by the RB for carrying out the matchmaking algorithm (to choose the “best” resource where to submit the job)
 - *Computing Resource*
 - Used to build expressions of Requirements and/or Rank attributes by the user
 - Have to be prefixed with “other.”
 - *Data and Storage resources*
 - Input data to process, SE where to store output data, protocols spoken by application when accessing Ses
- ◆ Note: GGF currently tries to standardise a Job Submission Description Language

JDL: relevant attributes

◆ **JobType**

- *Normal* (simple, sequential job), *Interactive*, *MPICH*, *Checkpointable*
- Or combination of them

◆ **Executable** (mandatory)

- The command name

◆ **Arguments** (optional)

- Job command line arguments

◆ **StdInput, StdOutput, StdError** (optional)

- Standard input/output/error of the job

◆ **Environment**

- List of environment settings

◆ **InputSandbox** (optional)

- List of files on the UI local disk needed by the job for running
- The listed files will automatically staged to the remote resource

◆ **OutputSandbox** (optional)

- List of files, generated by the job, which have to be retrieved

JDL: relevant attributes

◆ Requirements

- Job requirements on computing resources
- Specified using attributes of resources published in the Information Service
- If not specified, default value defined in UI configuration file is considered
 - Default: *other.GlueCEStateStatus* == "Production" (the resource has to be able to accept jobs and dispatch them on WNs)

◆ Rank

- Expresses preference (how to rank resources that have already met the Requirements expression)
- Specified using attributes of resources published in the Information Service
- If not specified, default value defined in the UI configuration file is considered
 - Default: - *other.GlueCEStateEstimatedResponseTime* (the lowest estimated traversal time)
 - Default: *other.GlueCEStateFreeCPUs* (the highest number of free CPUs) for parallel jobs (see later)

◆ GLUE Schema is used

JDL: relevant attributes

◆ **InputData**

- Refers to data used as input by the job: these data are published in the Replica Location Service (RLS) and stored in the SEs
- LFNs and/or GUIDs

◆ **DataAccessProtocol** (mandatory if InputData has been specified)

- The protocol or the list of protocols which the application is able to speak with for accessing *InputData* on a given SE

◆ **OutputSE**

- The Uniform Resource Identifier of the output SE
- RB uses it to choose a CE that is compatible with the job and is close to SE

Example of JDL File

```
[  
JobType="Normal";  
  
Executable = "gridTest";  
  
StdError = "stderr.log";  
  
StdOutput = "stdout.log";  
  
InputSandbox = {"home/joda/test/gridTest"};  
OutputSandbox = {"stderr.log", "stdout.log"};  
  
InputData = {"lfn:green", "guid:red"};  
  
DataAccessProtocol = "gridftp";  
  
Requirements = other.GlueHostOperatingSystemNameOpSys == "LINUX"  
                && other.GlueCEStateFreeCPUs>=4;  
  
Rank = other.GlueCEPolicyMaxCPUtime;  
]
```

General concepts of Grid Workload Management Systems

EGEE-0 Workload Management Systems

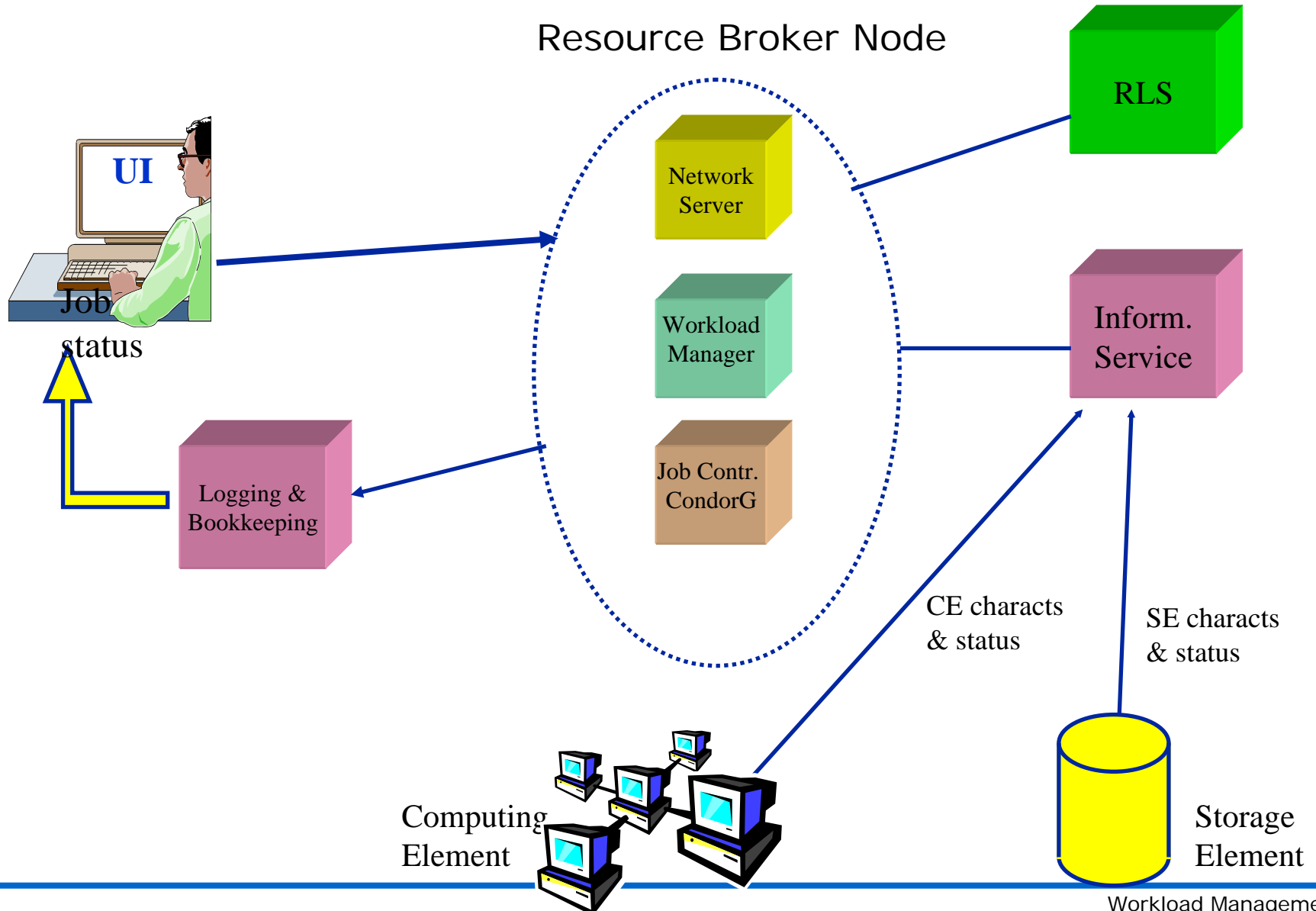
Job Preparation

Architecture / Job submission and status monitoring

Matchmaking

Different job types

Simplified Architecture Overview



Job Submission

```
edg-job-submit [-r <res_id>] [-c <config file>]  
[-vo <VO>] [-o <output file>] <job.jdl>
```

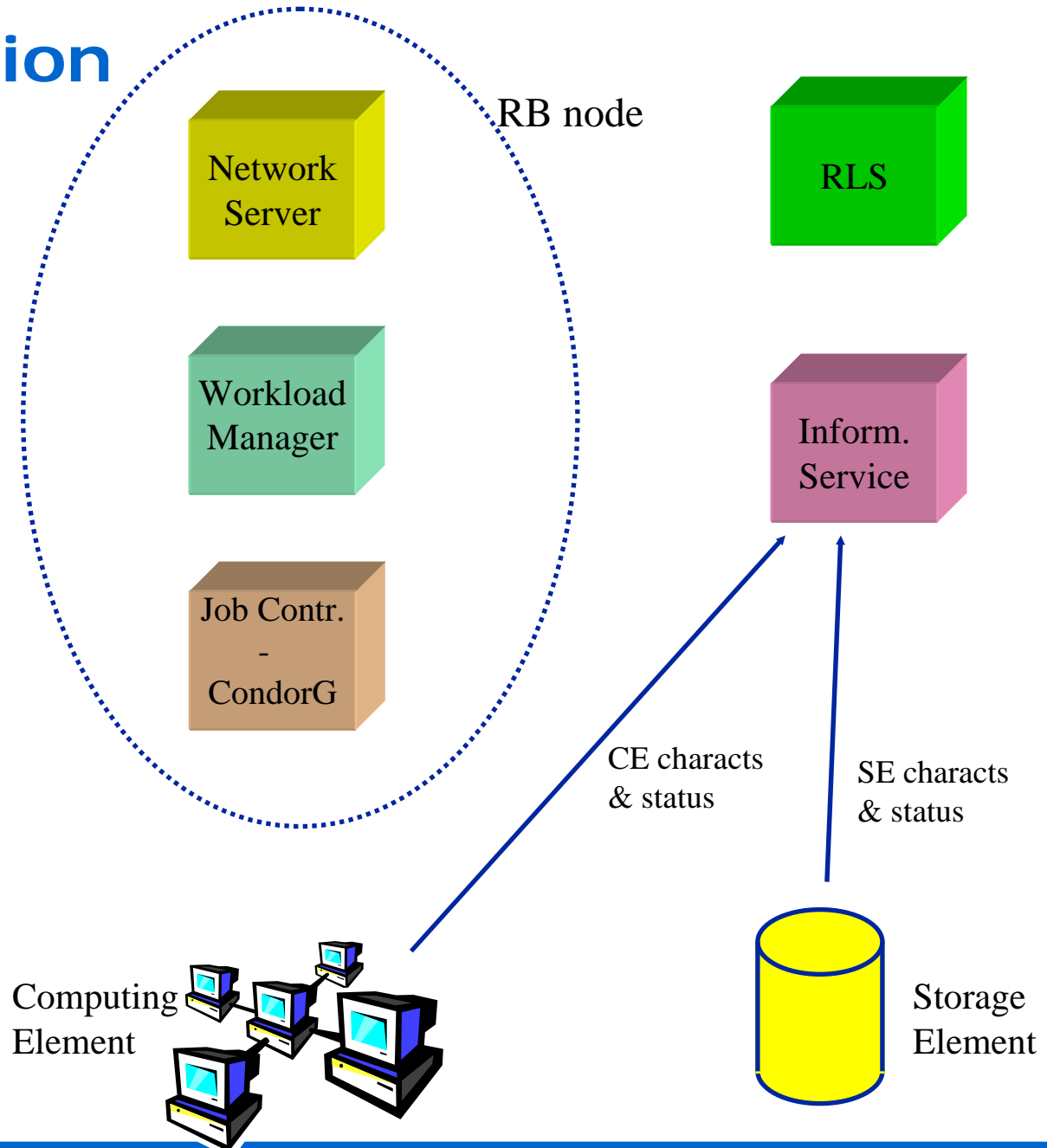
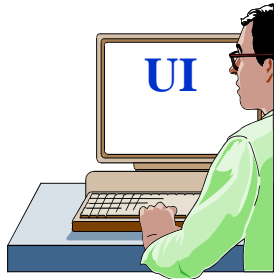
- r the job is submitted directly to the computing element identified by *<res_id>*
- c the configuration file *<config file>* is pointed by the UI instead of the standard configuration file
- vo the Virtual Organization (if user is not happy with the one specified in the UI configuration file)
- o the generated `edg_jobId` is written in the *<output file>*

Useful for other commands, e.g.:

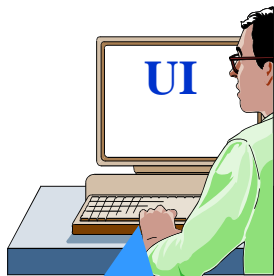
```
edg-job-status -i <input file> (or edg_jobId)
```

- i the status information about `edg_jobId` contained in the *<input file>* are displayed

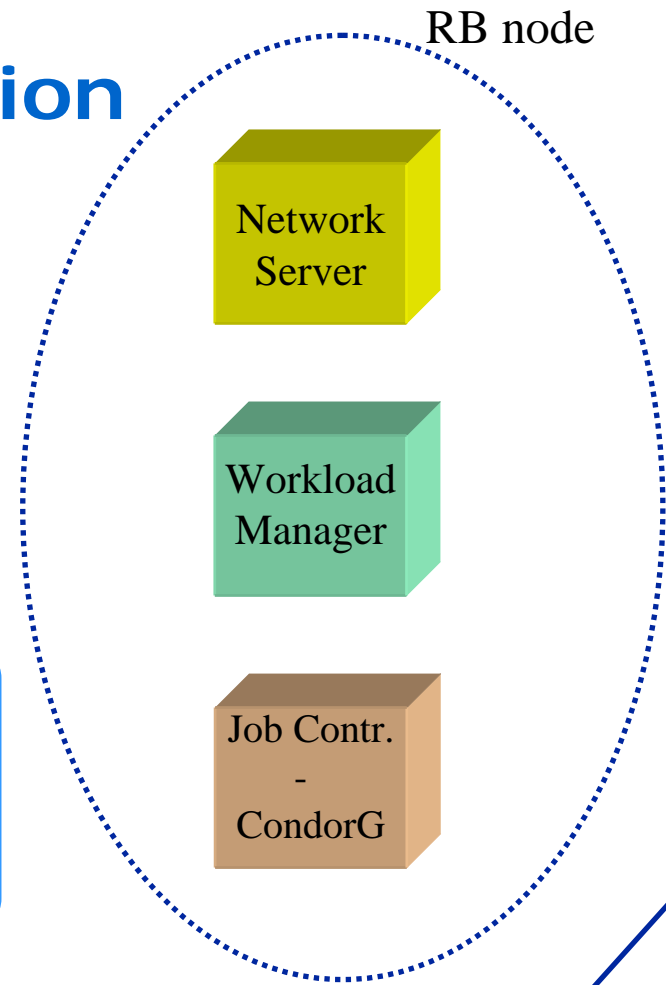
Job submission



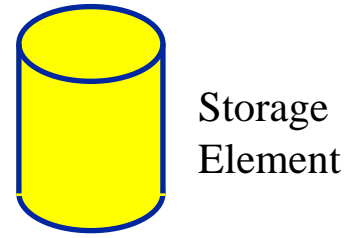
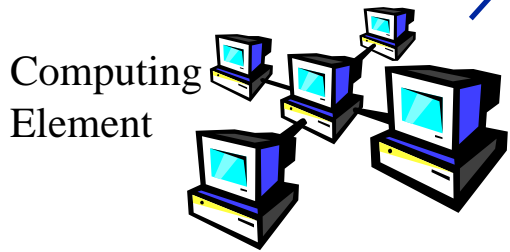
Job submission



UI: allows users to access the functionalities of the WMS (via command line, GUI, C++ and Java APIs)

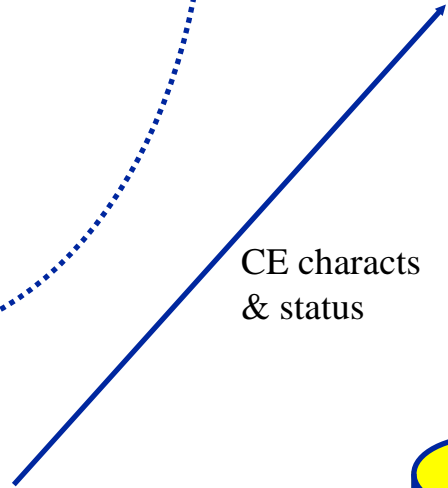


Job Status
submitted

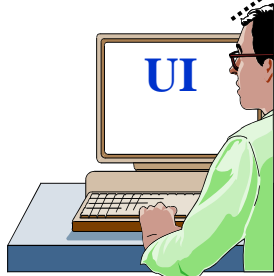


CE characts & status

SE characts & status



Job subm



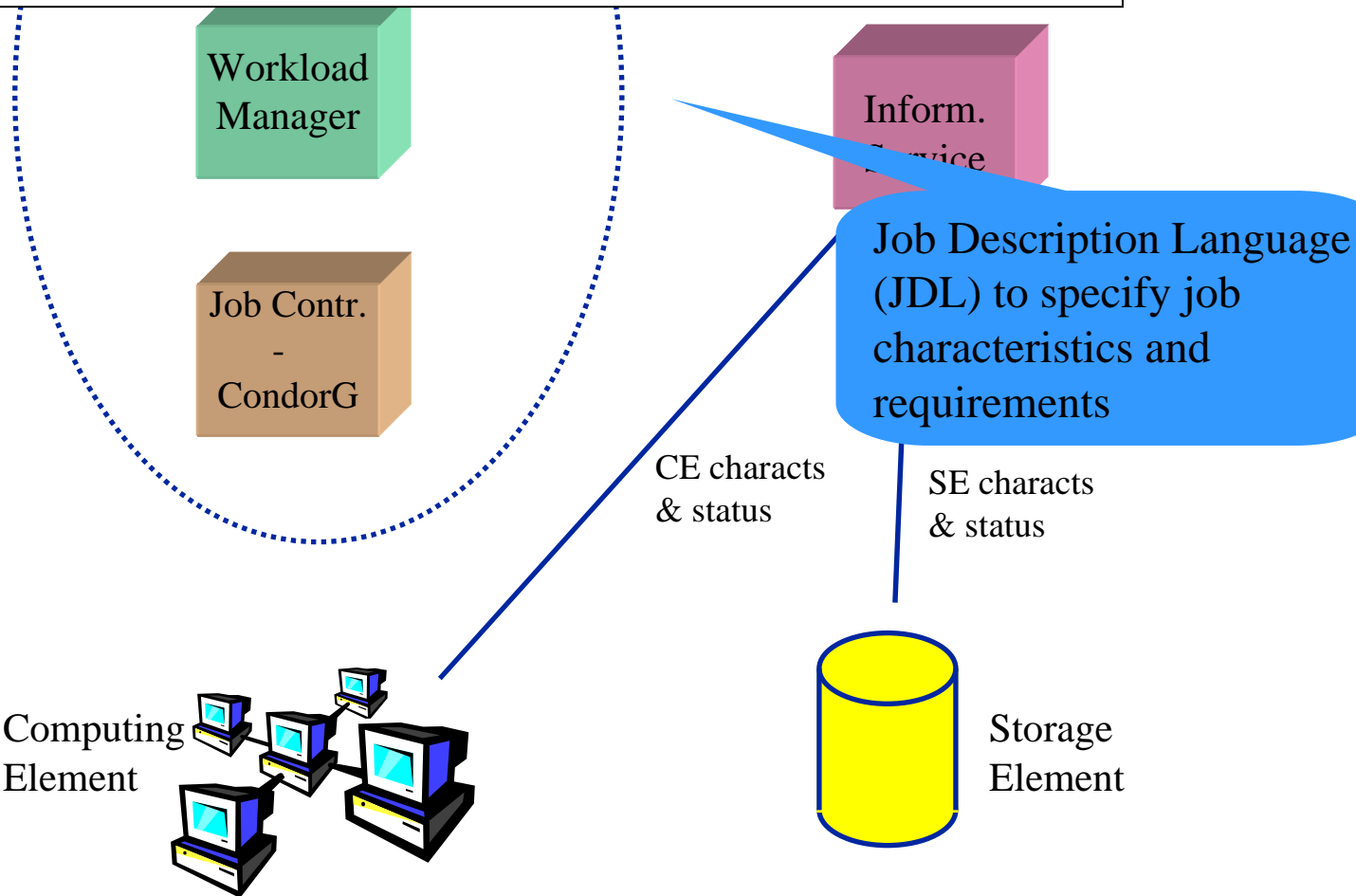
```
edg-job-submit myjob.jdl
```

```
Myjob.jdl
```

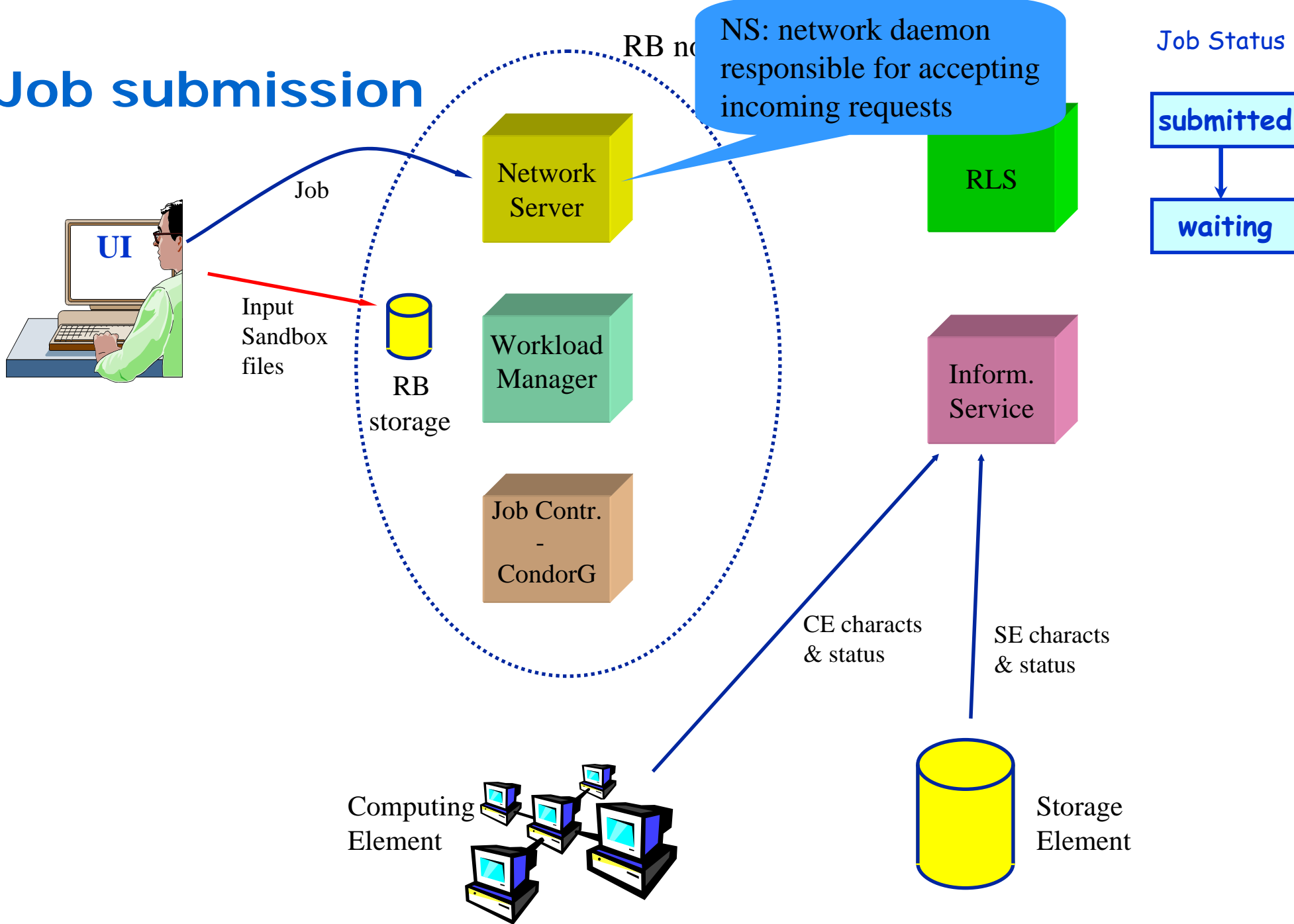
```
JobType = "Normal";  
Executable = "$(CMS)/exe/sum.exe";  
InputSandbox = {"/home/user/WP1testC", "/home/file*", "/home/user/DATA/*"};  
OutputSandbox = {"sim.err", "test.out", "sim.log"};  
Requirements = other.GlueHostOperatingSystemName == "linux" &&  
other.GlueHostOperatingSystemRelease == "Red Hat 7.3" &&  
other.GlueCEPolicyMaxWallClockTime > 10000;  
Rank = other.GlueCEStateFreeCPUs;
```

Job
Status

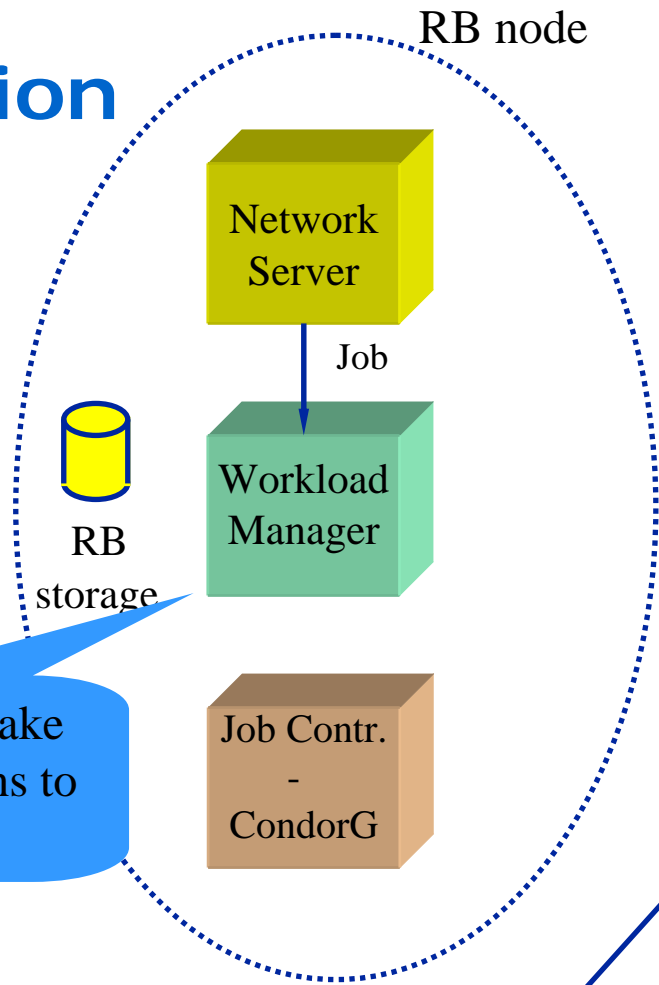
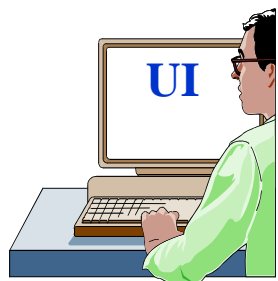
submitted



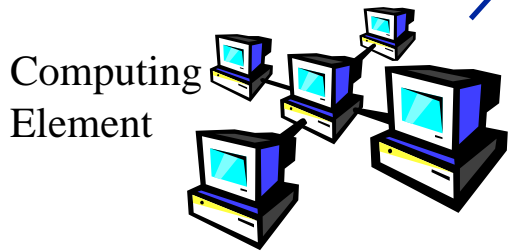
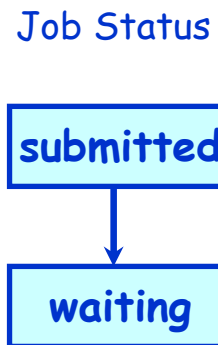
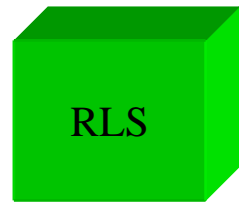
Job submission



Job submission

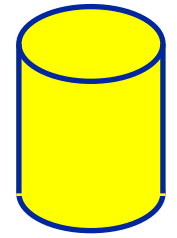


WM: responsible to take the appropriate actions to satisfy the request



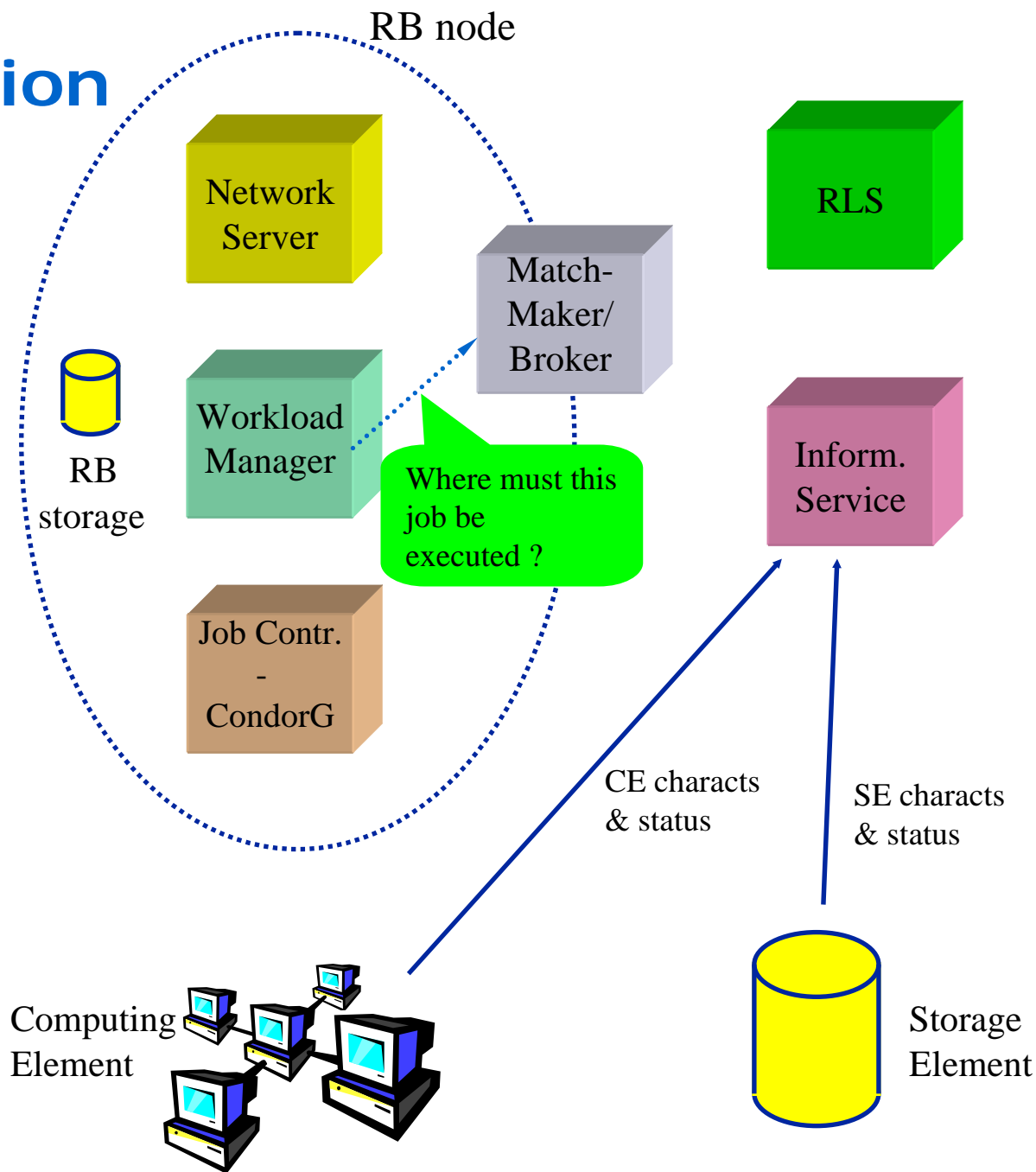
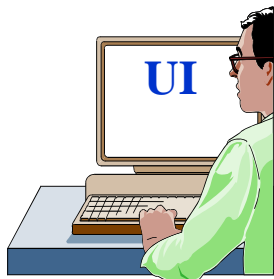
CE characts & status

SE characts & status



Storage Element

Job submission

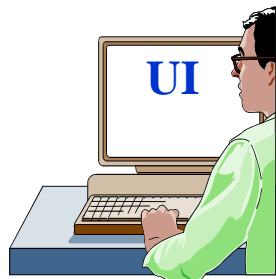


Job Status

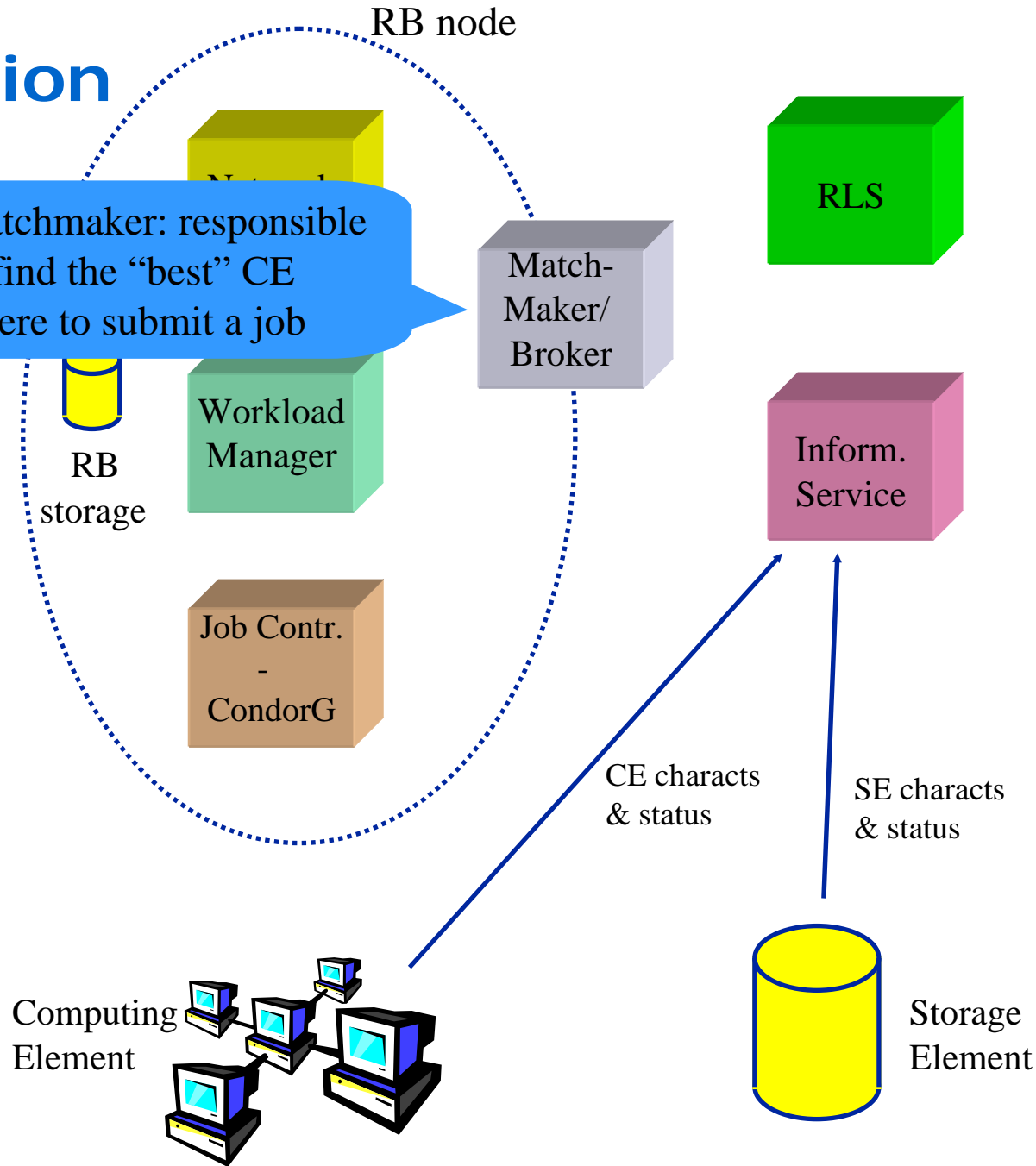
submitted

waiting

Job submission



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

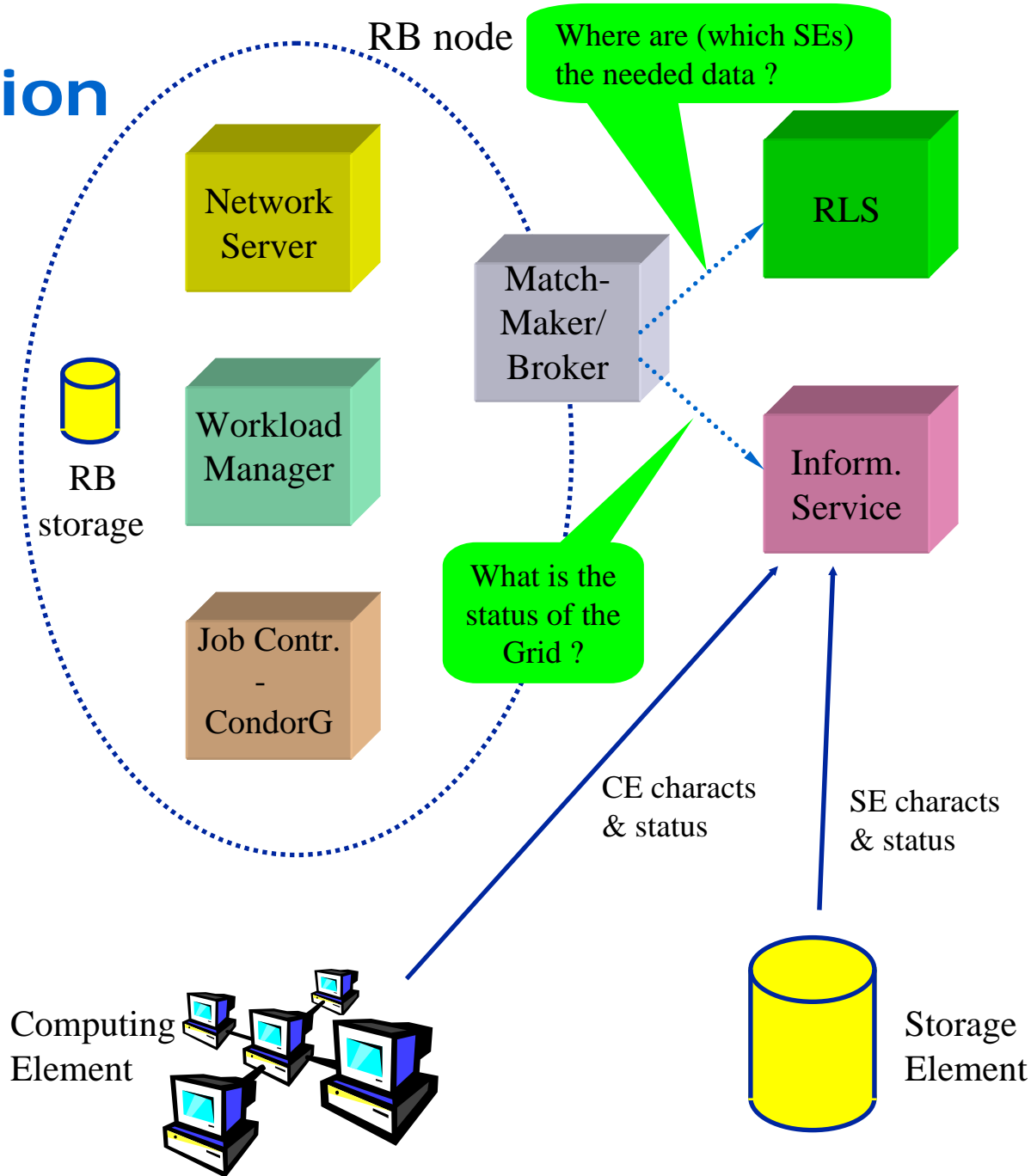
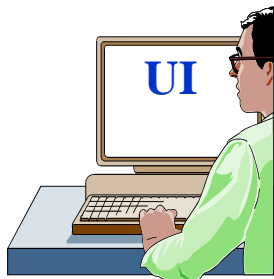


Job Status

submitted

waiting

Job submission

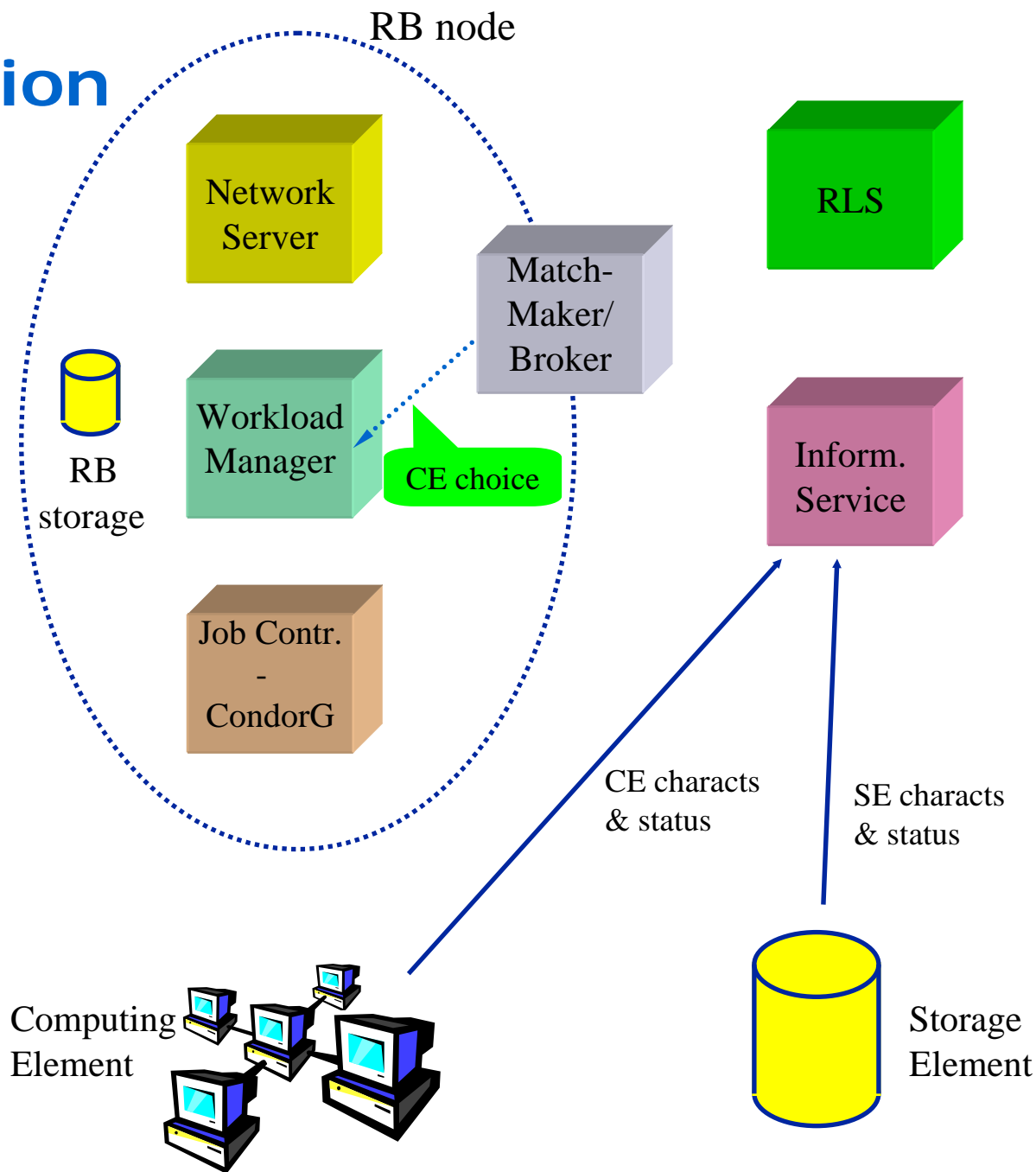
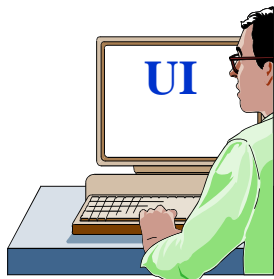


Job Status

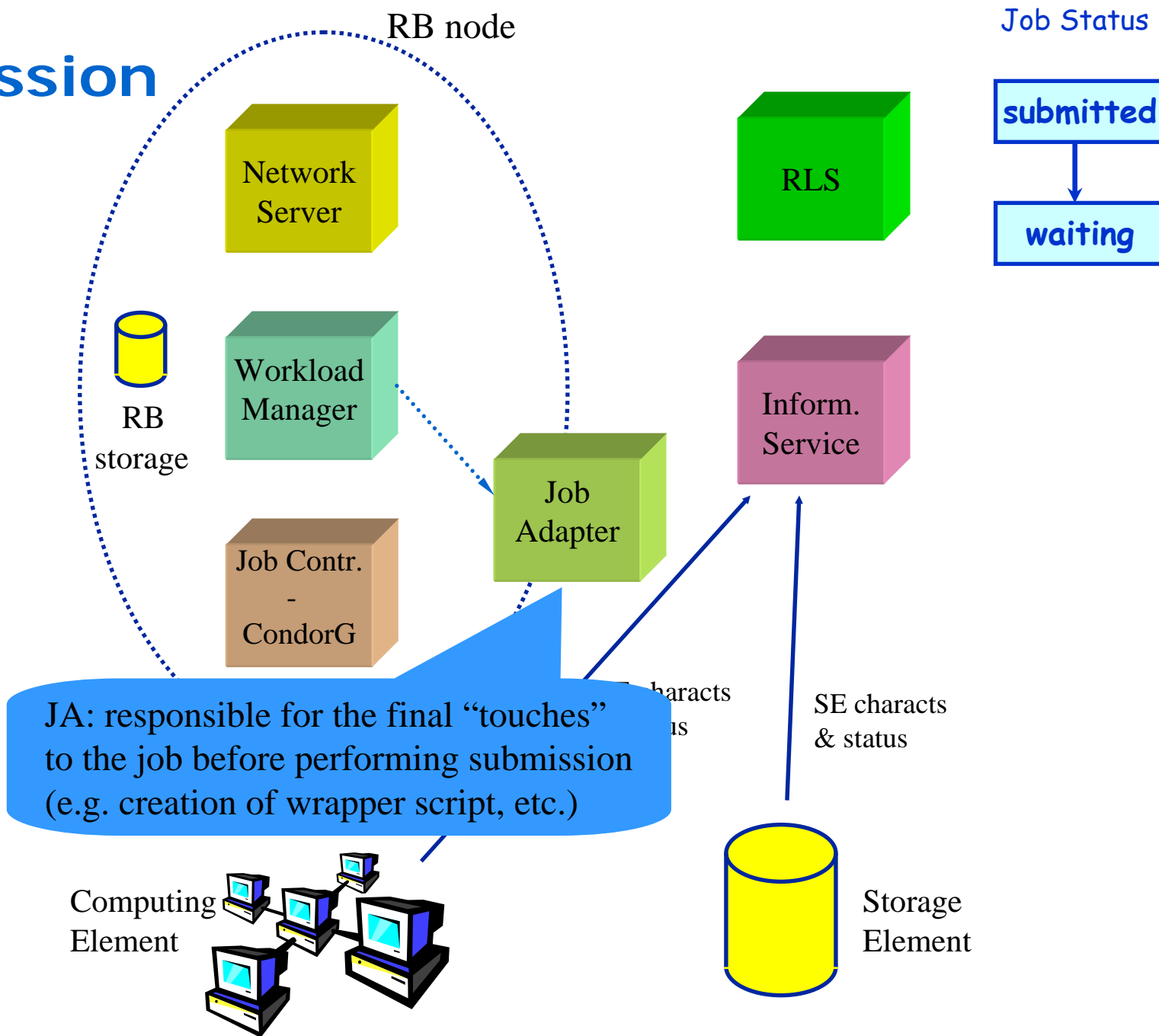
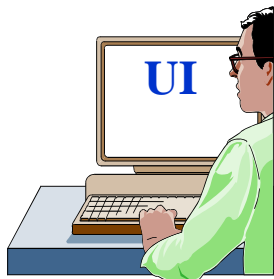
submitted

waiting

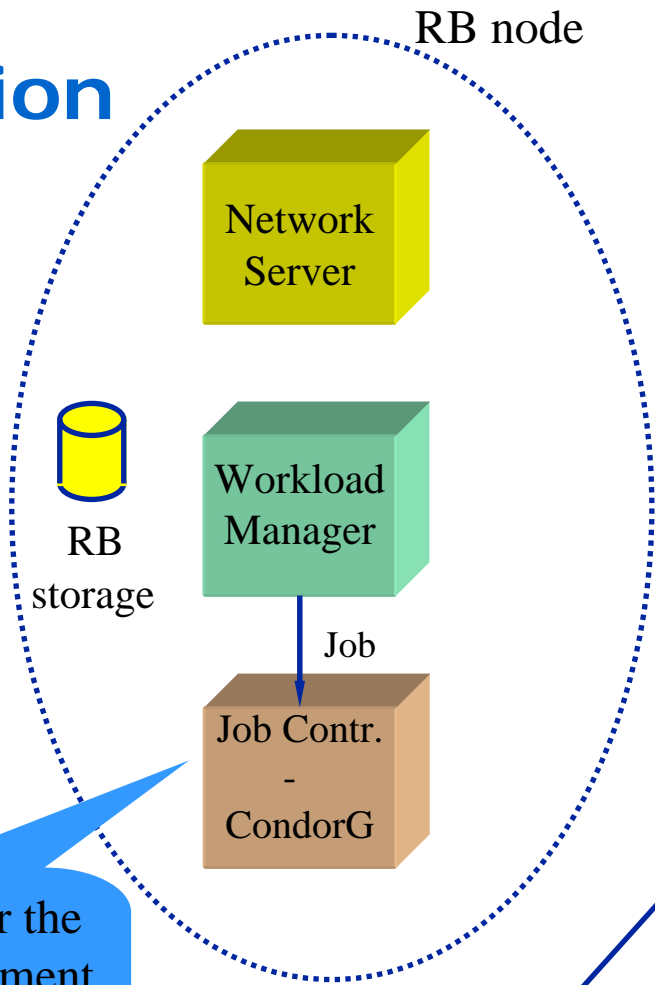
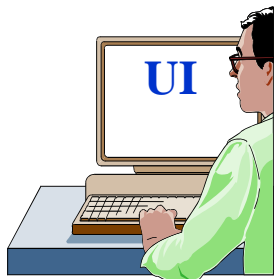
Job submission



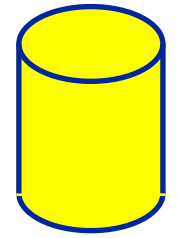
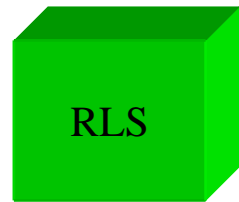
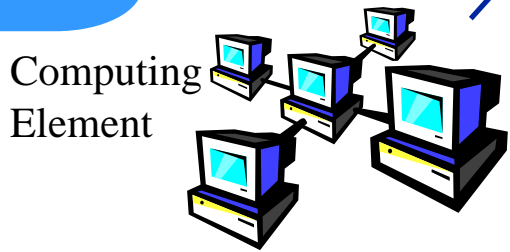
Job submission



Job submission



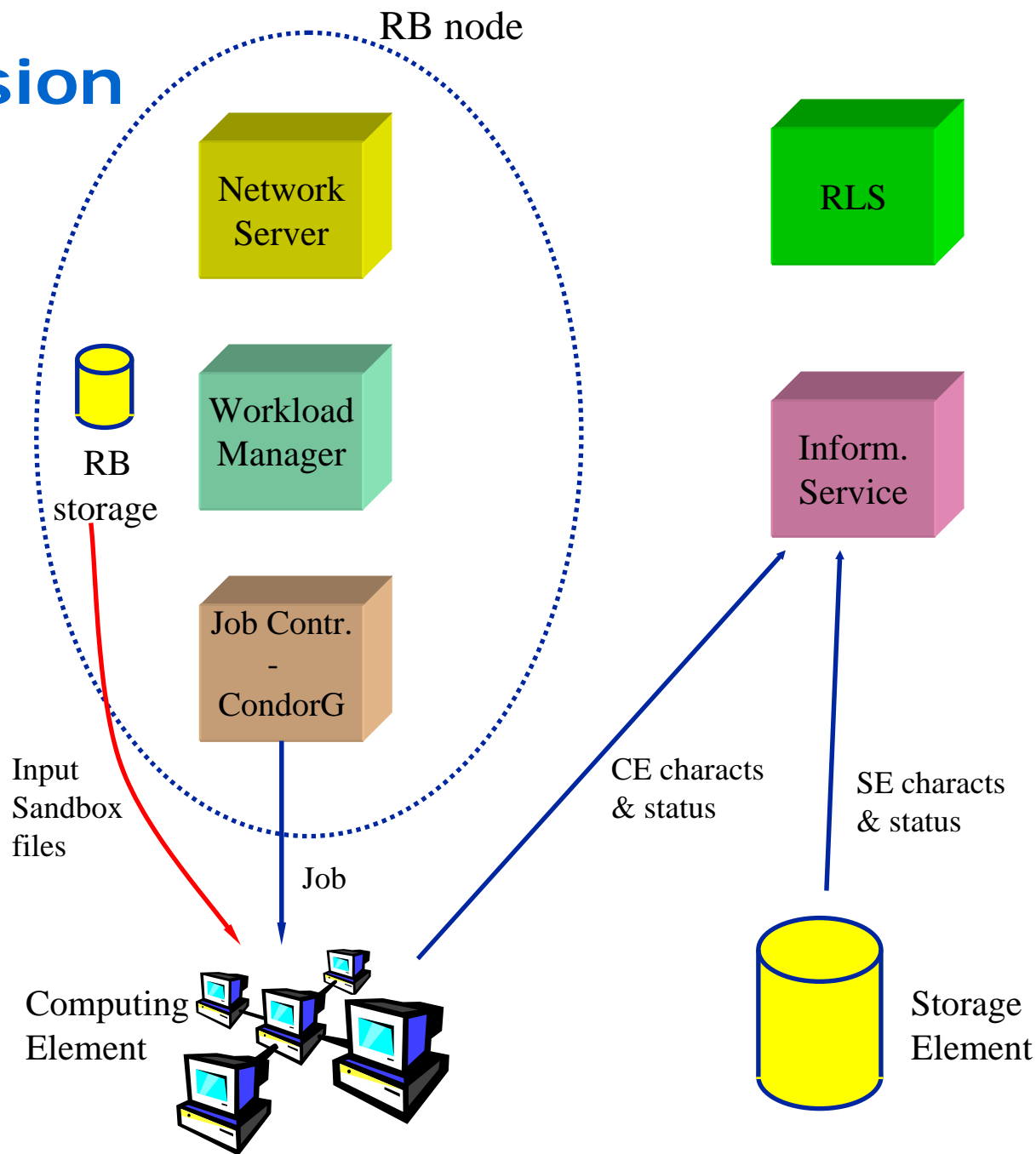
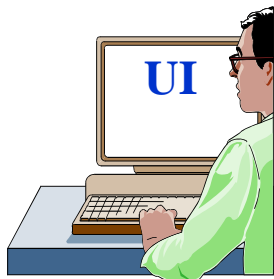
JC: responsible for the actual job management operations (done via CondorG)



CE characts & status

SE characts & status

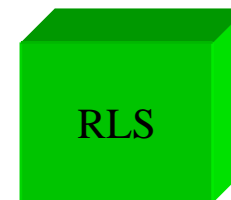
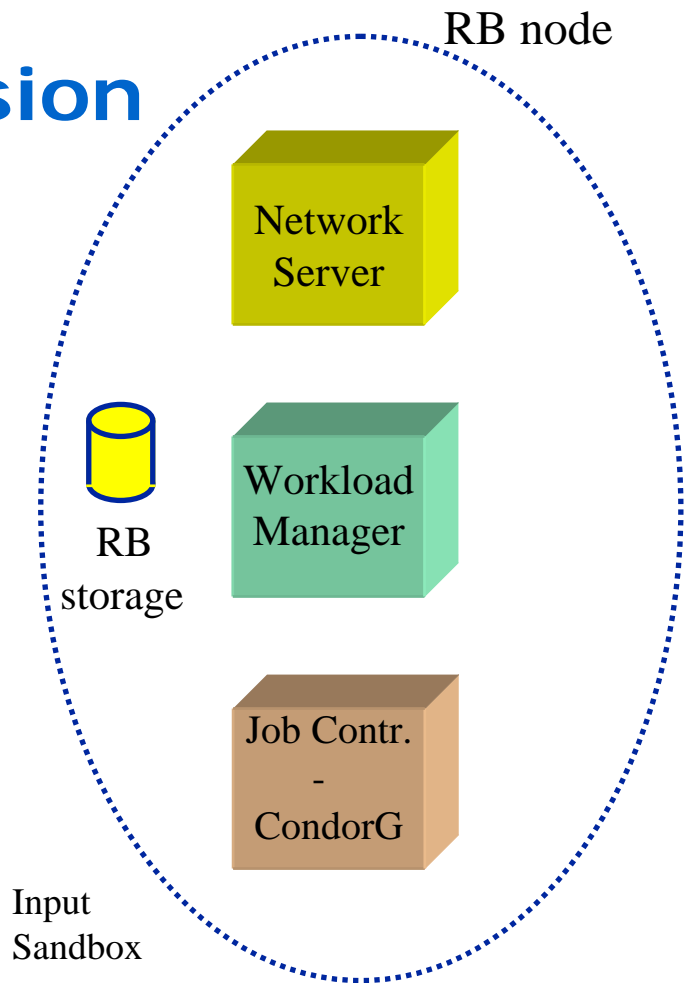
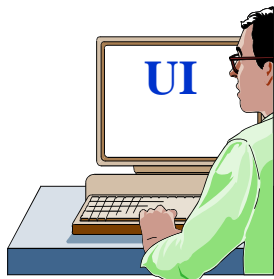
Job submission



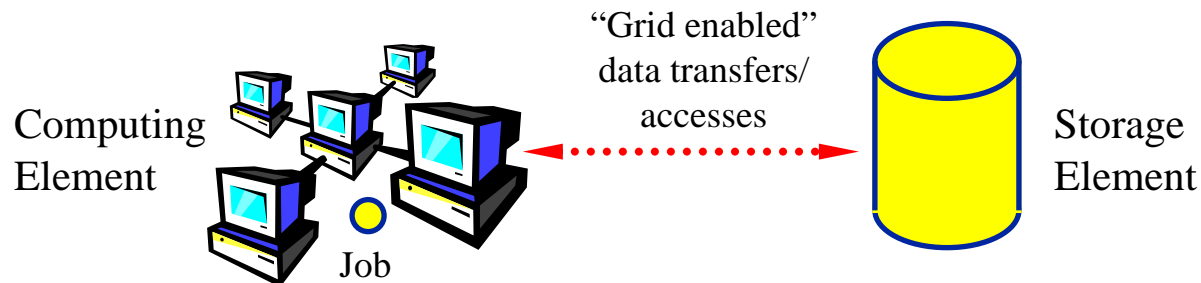
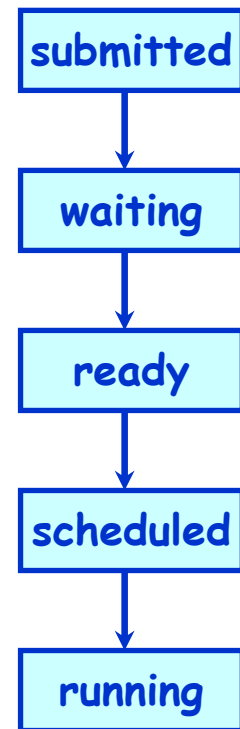
Job Status



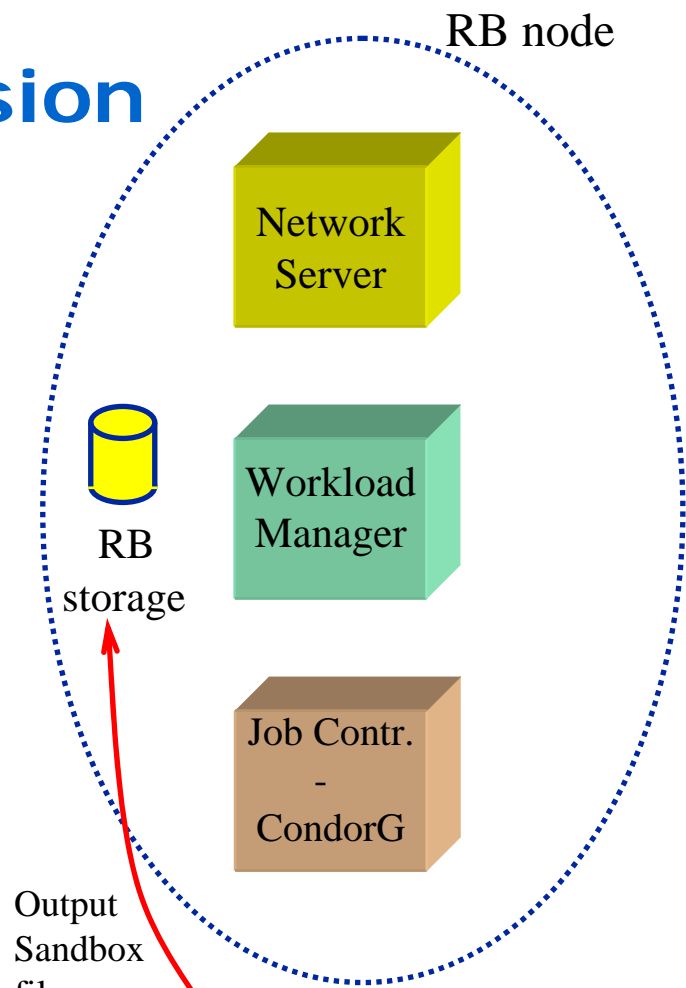
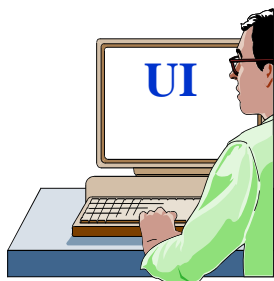
Job submission



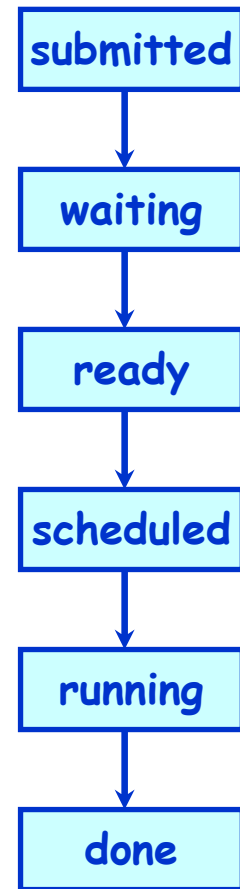
Job Status



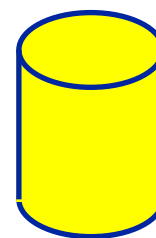
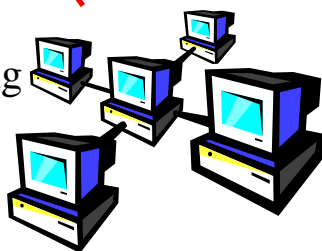
Job submission



Job Status



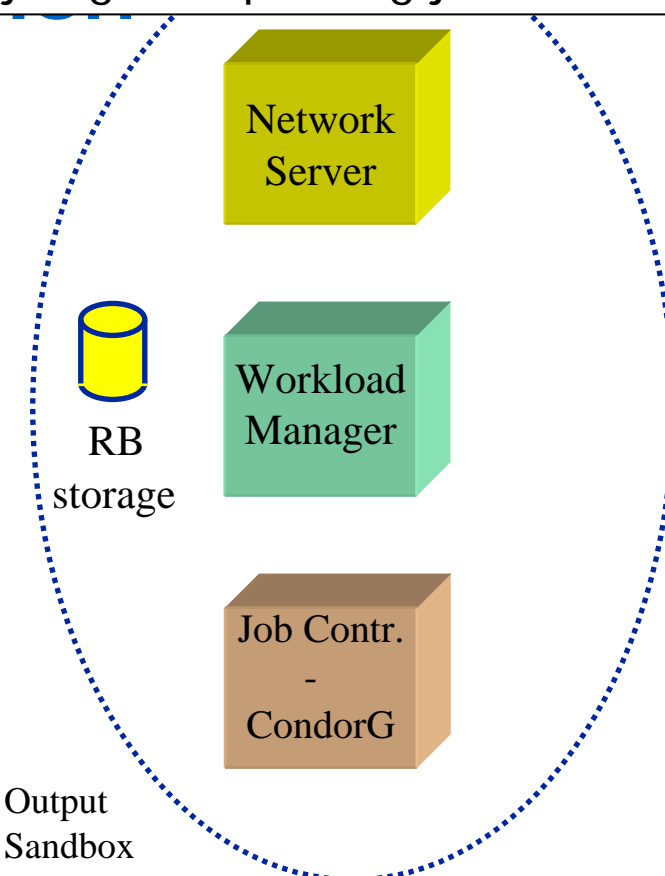
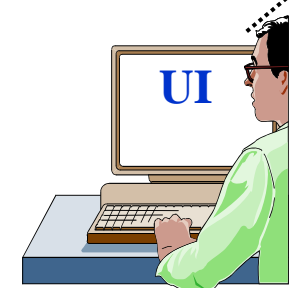
Computing Element



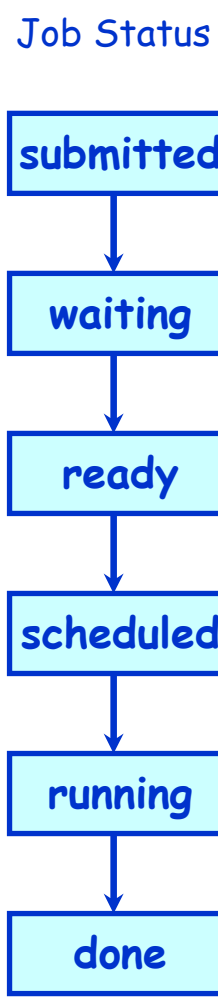
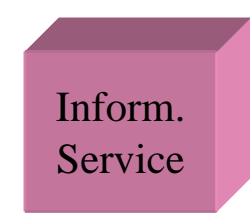
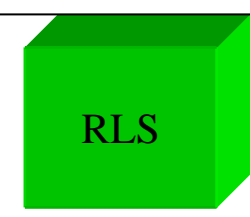
Storage Element

Job submission

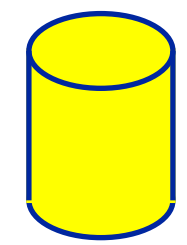
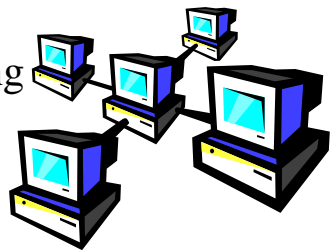
```
edg-job-get-output <dg-job-id>
```



RB node

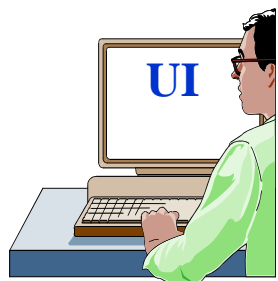


Computing Element



Storage Element

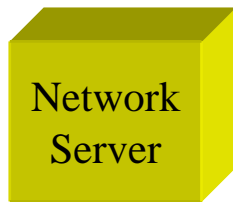
Job submission



Output
Sandbox
files



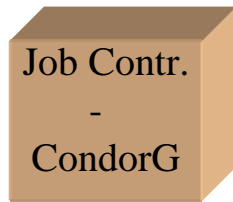
RB
storage



Network
Server



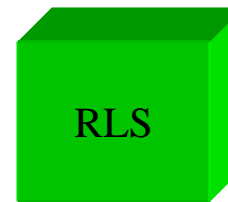
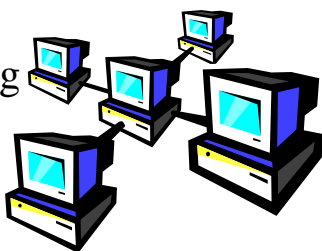
Workload
Manager



Job Contr.
-
CondorG

RB node

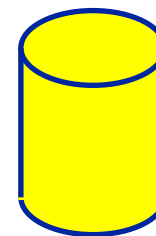
Computing
Element



RLS



Inform.
Service



Storage
Element

Job Status

submitted

waiting

ready

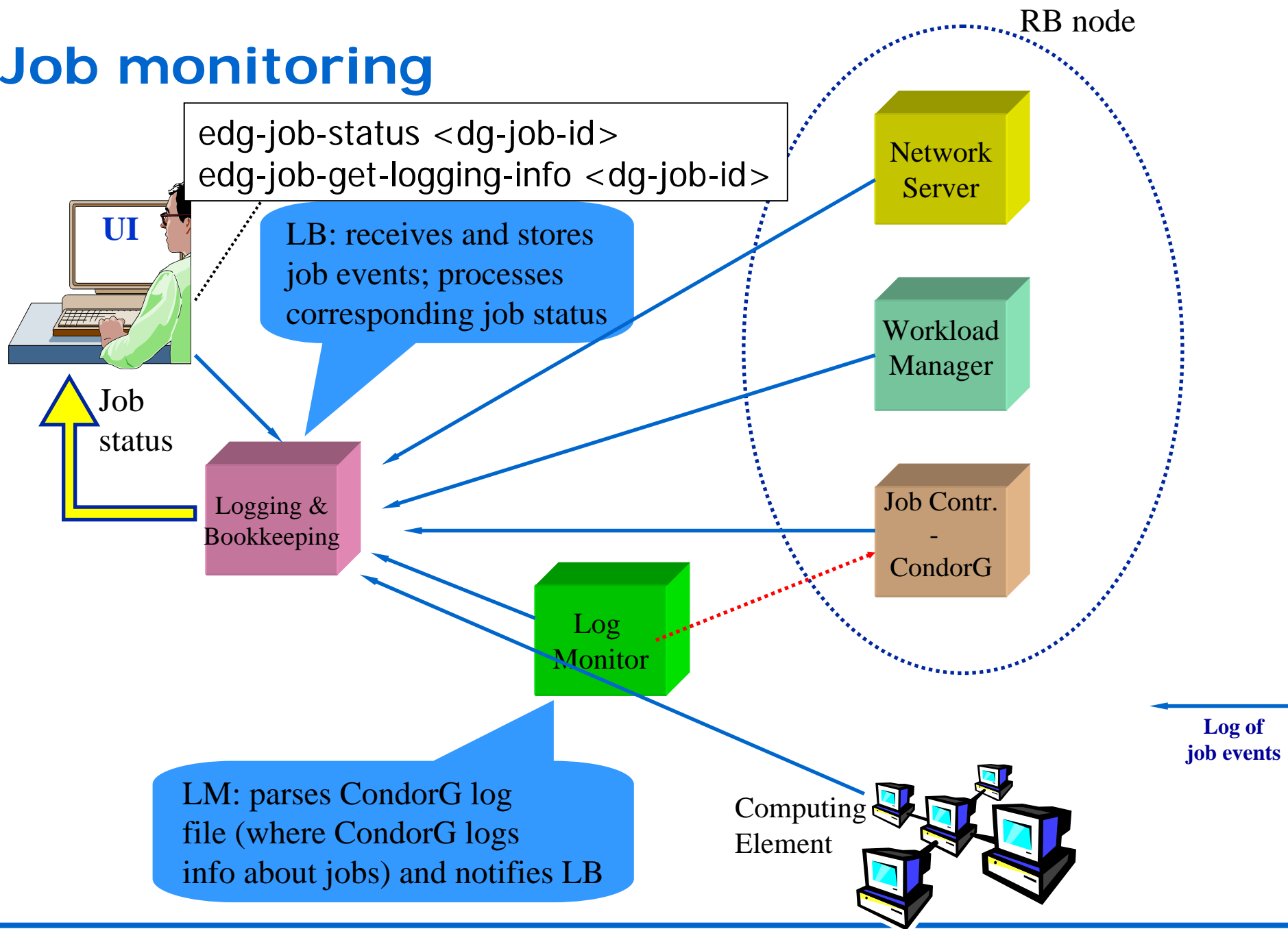
scheduled

running

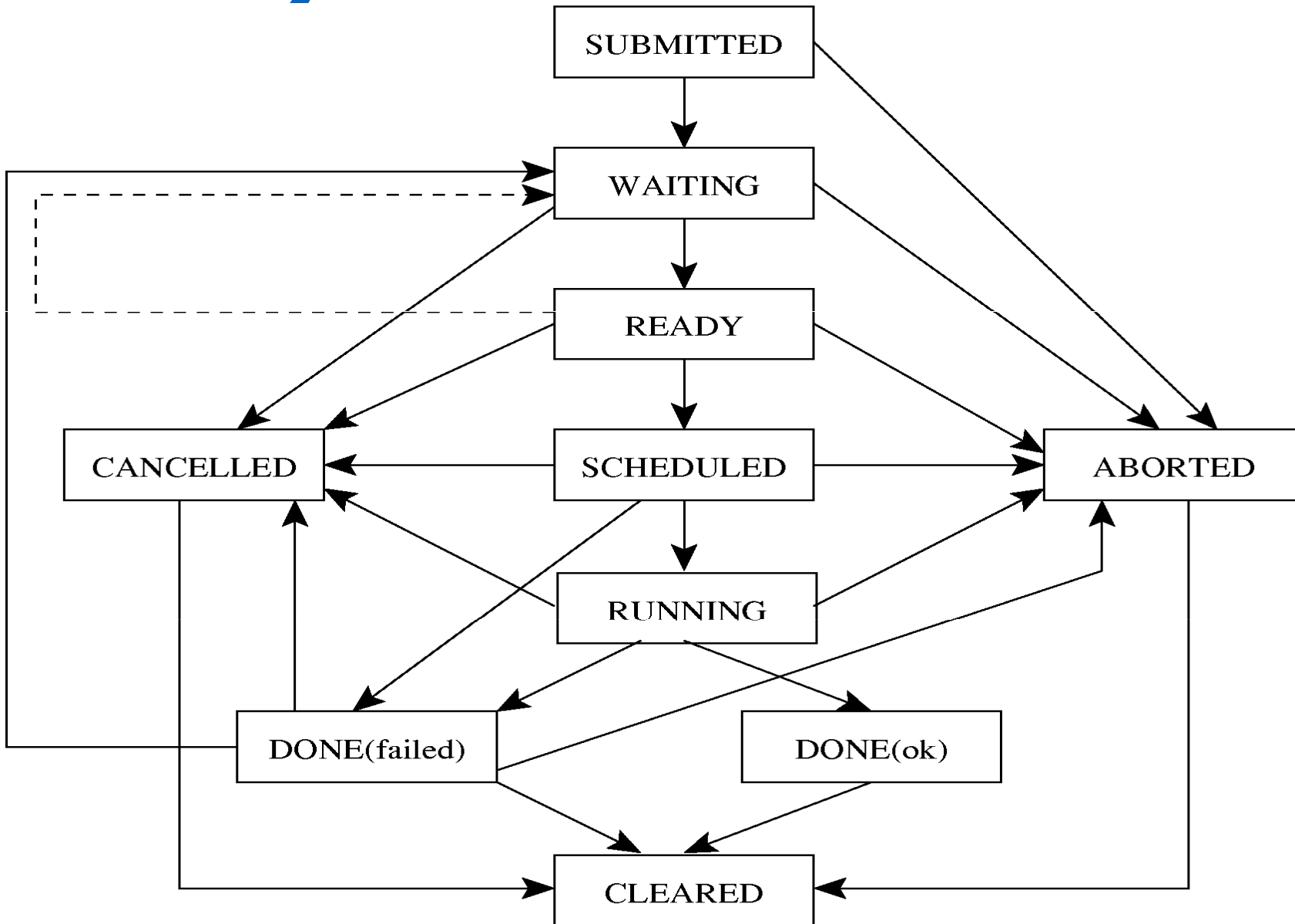
done

cleared

Job monitoring



Possible job states



Job resubmission

- ◆ If something goes wrong, the WMS tries to **reschedule and resubmit** the job (possibly on a different resource satisfying all the requirements)
- ◆ Maximum number of resubmissions: $\min(\text{RetryCount}, \text{MaxRetryCount})$
 - **RetryCount**: JDL attribute
 - **MaxRetryCount**: attribute in the "RB" configuration file
- ◆ E.g., to disable job resubmission for a particular job: *RetryCount=0*; in the JDL file

Proxy renewal

◆ Why?

- To avoid job failure because it outlived the validity of the initial proxy, avoiding considering long term user proxies

◆ Solution

- Short term proxies created as usual in the UI machine
 - *grid-proxy-init -hours <hours>*
- User registers proxy into a MyProxy server:
 - *myproxy-init -s <server> [-t <cred> -c <proxy>]*
 - *server* is the server address (e.g. lxshare0375.cern.ch)
 - *cred* is the number of hours the proxy should be valid on the server
 - *proxy* is the number of hours renewed proxies should be valid
- User specifies the MyProxy server in the JDL to enable proxy renewal:
 - `MyProxyServer=myproxy.host.name;`
- The Proxy is automatically renewed by WMS without user intervention for all the job life

Other (most relevant) UI commands

◆ **edg-job-list-match**

- Lists resources matching a job description
- Performs the matchmaking without submitting the job

◆ **edg-job-cancel**

- Cancels a given job

◆ **edg-job-status**

- Displays the status of the job

◆ **edg-job-get-output**

- Returns the job-output (the OutputSandbox files) to the user

◆ **edg-job-get-logging-info**

- Displays logging information about submitted jobs (all the events “pushed” by the various components of the WMS)
- Very useful for debug purposes

General concepts of Grid Workload Management Systems

EGEE-0 Workload Management Systems

Job Preparation

Architecture / Job submission and status monitoring

Matchmaking

Different job types

WMS Matchmaking

- ◆ The RB (Matchmaker) has to find the best suitable computing resource (CE) where the job will be executed
- ◆ It interacts with Data Management Service and Information Services
 - They supply RB with all the information required for the resolution of the matches
- ◆ The CE chosen by RB has to match the job requirements (e.g. runtime environment, data access requirements, and so on)
- ◆ If *FuzzyRank=False* (default):
 - If 2 or more CEs satisfy all the requirements, the one with the best Rank is chosen
 - If there are two or more CEs with the same best rank, the choice is done in a random way among them
- ◆ If *FuzzyRank=True* in the JDL:
 - Fuzziness in CE choice: the CE with highest rank has the highest probability to be chosen

WMS matchmaking scenarios

◆ Possible scenarios for matchmaking:

1. Direct job submission

- `edg-job-submit -r <CEId>`
- Corresponds to job submission with Globus clients (`globus-job-submit`)

2. Job submission with computational requirements only

- Nor `InputData` nor `OutputSE` specified in the JDL

3. Job submission with data access requirements

- `InputData` and/or `OutputSE` specified in the JDL

4. Matchmaking with `getAccessCost`

Direct job submission

edg-job-submit -r CEId

- ◆ Job is simply submitted on the given CE
- ◆ RB does not perform any matchmaking algorithm
- ◆ Information services not queried at all

Job submission with only comput. reqs

- ◆ Nor `InputData` nor `OutputSE` specified in the JDL

- ◆ Matchmaking algorithm:
 - Requirements check
 - RB contacts the IS to check which CEs satisfy all the requirements
 - This includes also authorization check (where is the user allowed to submit jobs ?)
 - Suitable resources directly queried (GRISes queried) to evaluate Rank expression (which usually refers to dynamic values)
 - If more than one CE satisfies the job requirements, the CE with the best rank is chosen by the RB (or has the highest probability to be chosen, if `Fuzzyrank` enabled)

Job submission with data access reqs

- ◆ `InputData` and/or `OutputSE` specified in the JDL
- ◆ RB strategy: submit jobs close to data
- ◆ Matchmaking algorithm:
 - Requirements check as in the previous case
 - CE chosen among the suitable ones (the CEs which passed the requirements check) and where most of the needed files are “close” to it (where most of the needed files are stored on SEs close to the considered CE)

Matchmaking with GetAccessCost

- ◆ Can be used when `InputData` has been specified in the JDL
- ◆ Used when `Rank = other.DataAccessCost` has been specified in the JDL
- ◆ Matchmaking algorithm:
 - Requirements check as in the previous case
 - The CE is chosen by the 'getAccessCost' method provided by data Management Services among the suitable CEs (the CEs which passed the requirements check), taking into account data location and network information

Example of job submission

- ◆ User logs in on the UI
- ◆ User issues a *grid-proxy-init* and enters his certificate's password, getting a valid Globus proxy
- ◆ User sets up his or her JDL file
- ◆ Example of Hello World JDL file :

```
[  
  Executable = "/bin/echo";  
  Arguments = "Hello World";  
  StdOutput = "Message.txt";  
  StdError = "stderr.log";  
  OutputSandbox = { "Message.txt", "stderr.log" };  
]
```

Example of job submission

- ◆ User issues a: *edg-job-submit HelloWorld.jdl*
and gets back from the system a unique Job Identifier (JobId)

- ◆ User issues a: *edg-job-status JobId*
to get logging information about the current status of his Job

- ◆ When the "Output" status is reached, the user can issue a
edg-job-get-output JobId
and the system returns the name of the temporary directory where the job output can be found on the UI machine.

Example of job submission

```
$ edg-job-submit HelloWorld.jdl
```

```
*****
```

JOB SUBMIT OUTCOME

The job has been successfully submitted to the Network Server.

Use `edg-job-status` command to check job current status. Your job identifier (`edg_jobId`) is:

```
- https://lxshare0403.cern.ch:9000/KoBA-IgxZyVpLKhANfrhHw
```

```
*****
```



JobId

Example of job submission

```
$ edg-job-status https://lxshare0403.cern.ch:9000/KoBA-IgxZyVpLKhANfrhHw
```

```
*****
```

BOOKKEEPING INFORMATION:

```
Printing status info for the Job : https://lxshare0403.cern.ch:9000/KoBA-IgxZyVpLKhANfrhHw
```

```
Current Status:   Done (Success)
```

```
Exit code:       0
```

```
Status Reason:   Job terminated successfully
```

```
Destination:     lxshare0405.cern.ch:2119/jobmanager-pbs-infinite
```

```
reached on:     Wed Jun 18 12:06:10 2003
```

```
*****
```

Example of job submission

```
$ edg-job-get-output --dir Results https://lxshare0403.cern.ch:9000/KoBA-IgxZyVpLKhANfrhHw
```

```
*****
```

JOB GET OUTPUT OUTCOME

Output sandbox files for the job:

- <https://lxshare0403.cern.ch:9000/KoBA-IgxZyVpLKhANfrhHw>

have been successfully retrieved and stored in the directory:

</shift/lxshare072d/data01/UIhome/sgaravat/Results/KoBA-IgxZyVpLKhANfrhHw>

```
*****
```

```
$ more Results/KoBA-IgxZyVpLKhANfrhHw/Message.txt
```

```
Hello World
```

```
$ more Results/KoBA-IgxZyVpLKhANfrhHw/stderr.log
```

```
$
```


General concepts of Grid Workload Management Systems

EGEE-0 Workload Management Systems

Job Preparation

Architecture/ Job submission and status monitoring

Matchmaking

Different job types

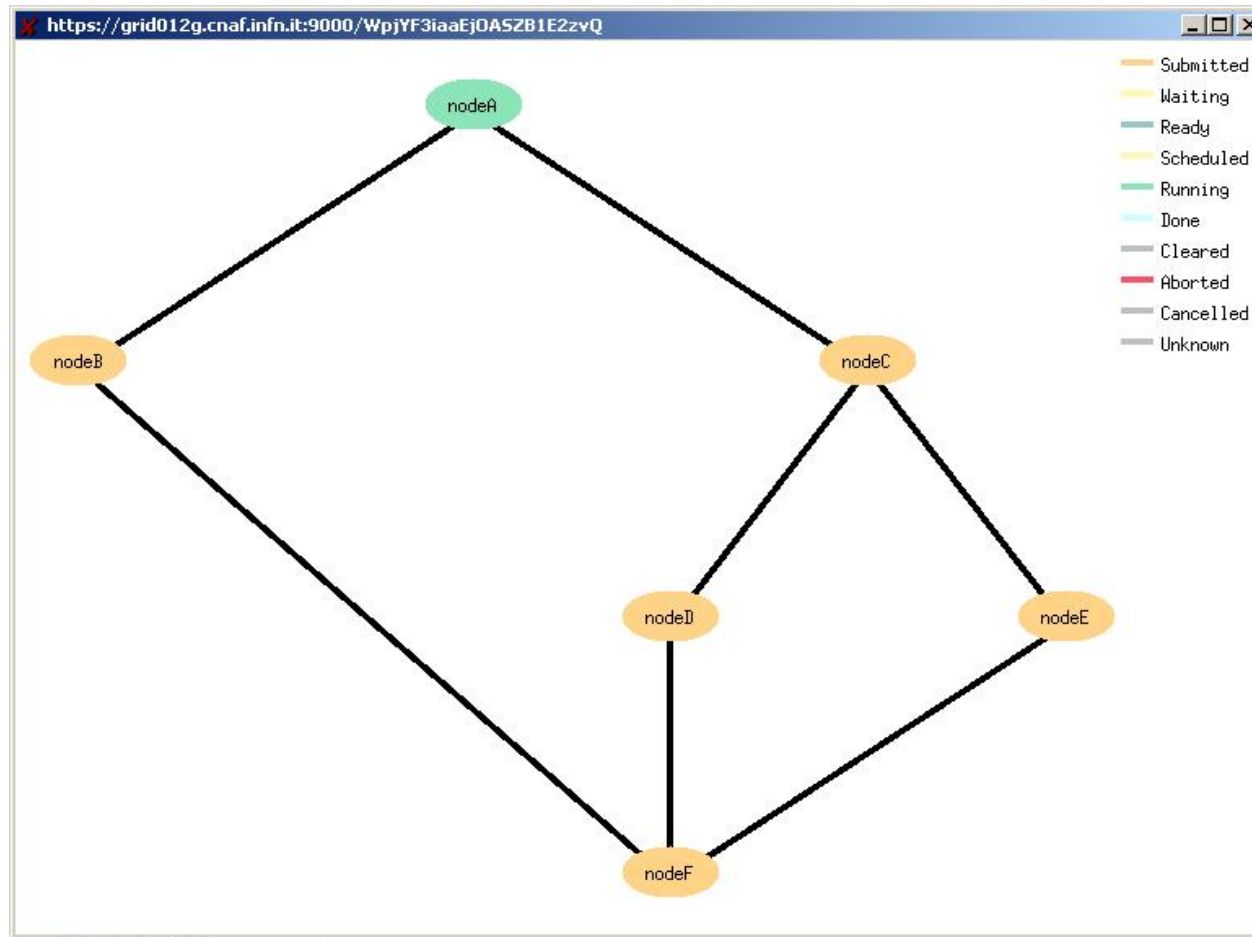
Interactive jobs

- ◆ Specified setting `JobType = "Interactive"` in JDL
- ◆ When an interactive job is executed, a window for the stdin, stdout, stderr streams is opened
 - Possibility to send the stdin to the job
 - Possibility to have the stderr and stdout of the job when it is running
- ◆ Possibility to start a window for the standard streams for a previously submitted interactive job with command `edg-job-attach`



Job Dependencies

Condor's **DAGman** allows for job dependencies



DAG = Direct Acyclic Graph

Job check-pointing

- ◆ Check-pointing: saving from time to time job state
 - Useful to prevent data loss, due to unexpected failures
 - Approach: provide users with a “trivial” logical job check-pointing service
 - User can save from time to time the state of the job (defined by the application)
 - A job can be restarted from an intermediate (i.e. “previously” saved) job state
- ◆ Different than “classical check-pointing (i.e. saving all the information related to a process: process’s data and stack segments, open files, etc.)”
 - Very difficult to apply (e.g. problems to save the state of open network connections)
 - Not necessary for many applications
- ◆ To submit a check-pointable job
 - Code must be instrumented (see next slides)
 - `JobType=Checkpointable` to be specified in JDL

Job check-pointing scenarios

◆ Scenario 1

- Job submitted to a CE
- When job runs it saves from time to time its state
- Job failure, due to a Grid problems (e.g. CE problem)
- Job resubmitted by the WMS possibly to a different CE
- Job restarts its computation from the last saved state
 - → No need to restart from the beginning
 - → The computation done till that moment is not lost

◆ Scenario 2

- Job failure, but not detected by the Grid middleware
- User can retrieve a saved state for the job (typically the last one)
 - *edg-job-get-chkpt -o <state> <edg-jobid>*
- User resubmits the job, specifying that the job must start from a specific (the retrieved one) initial state
 - *edg-job-submit -chkpt <state> <JDL file>*

Submission of parallel jobs

- ◆ Possibility to submit **MPI** jobs
- ◆ MPICH implementation supported
- ◆ Only parallel jobs inside a single CE can be submitted
- ◆ Submission of parallel jobs very similar to normal jobs
 - Just needed to specify in the JDL:
 - `JobType= "MPICH"`
 - `NodeNumber = n;`
 - The number (n) of requested CPUs
- ◆ Matchmaking
 - CE chosen by RB has to have MPICH sw installed, and at least n total CPUs
 - If there are two or more CEs satisfying all the requirements, the one with the highest number of free CPUs is chosen

Further information

- ◆ The EDG User's Guide

<http://marianne.in2p3.fr>

- ◆ EDG WP1 Web site

<http://www.infn.it/workload-grid>

In particular WMS User & Admin Guide and JDL docs

- ◆ Condor ClassAd

<https://www.cs.wisc.edu/condor/classad>

Abbreviations

- ◆ GRAM – Grid Resource Access Manager
- ◆ JDL – Job Description Language
- ◆ RB – Resource Broker
- ◆ WMS – Workload Management System
- ◆ UI – User Interface