

Simple, fault tolerant, lightweight grid computing approach for bag-of-tasks applications

Yiannis Georgiou, Nicolas Capit, Bruno Bzeznik and
Olivier Richard

3rd EGEE User Forum



Grenoble, FRANCE
Mescal Project

Plan

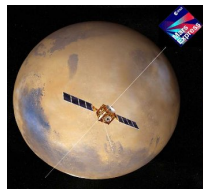
- 1 CIGRI: lightweight grid computing approach
- 2 Fault Treatment and Turnaround time Optimisations
- 3 Evaluating CIGRI upon Grid5000 platform
- 4 Conclusions

Plan

- 1 CIGRI: lightweight grid computing approach
 - Motivations and Related Work
 - CIGRI lightweight grid system: Principal Concepts
 - Global Architecture
- 2 Fault Treatment and Turnaround time Optimisations
- 3 Evaluating CIGRI upon Grid5000 platform
- 4 Conclusions

Motivations and Related Work

CIMENT project: **Mutualise the computing power** of Rhone-Alpes regional private laboratories **cluster resources** from different disciplines (environment, chemistry, physics, astronomy, medicine, ...) to effectuate **larger scale computations**



- Option of **Globus**: complicated, expensive (ex. Condor-G, Nimrod-G,...)
- Emergence of **desktop grid** systems and the idea of **cycle stealing** technologies provided the good bases for the CIGRI approach
- Alternative grid solutions: **Condor** (..low security,.. parallel applications), **OurGrid** (..high security, .."BoT" applications)

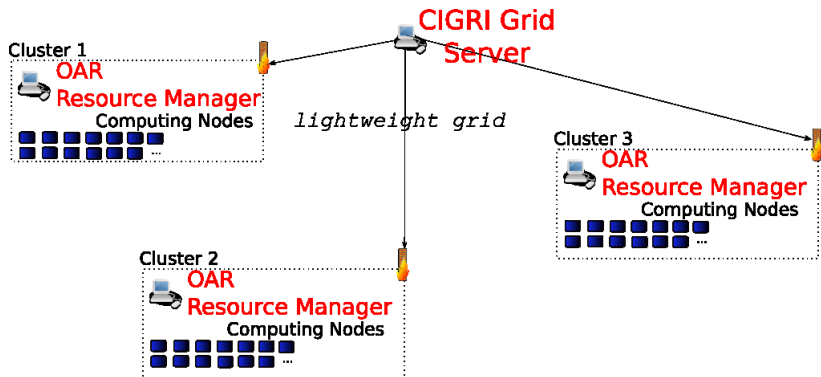
CIGRI approach for grid computing

- **CIGRI**: simpler, **lightweight approach** of grid platform (low security,..only "Bag-of-Tasks" applications (set of independent tasks))
- Using the method of aggregation of idle cluster resources
- Platform **focuses** on research and development of problems that come along with the execution of tasks (scheduling, fault tolerance,...)
- Choices of simple solutions on important classic grid issues like security, authentication mechanisms , resource location...



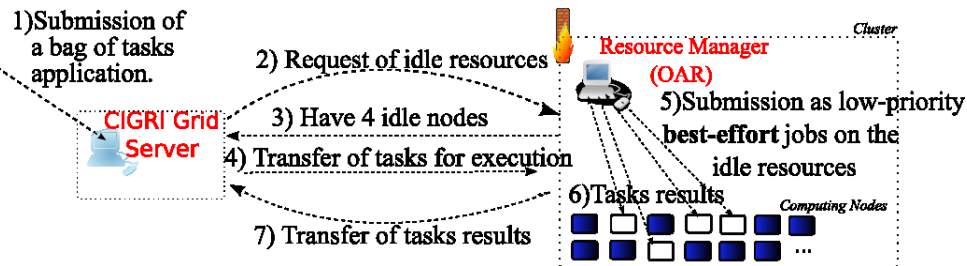
CIGRI: Lightweight grid

- Homogenisation of services and administration procedures between clusters



CIGRI: Cluster utilisation policy

- Works **discreetly** along with the interconnected clusters: **No specific** CIGRI software installed on clusters
- **Besteffort jobs**: job type provided by the cluster resource management system (introduced in OAR)
 - Lowest priority jobs submitted only if there is a free resource



- Killed when local cluster job requests the resource

CIGRI Global Architecture

- **High level components:** MySQL Database, Perl programming language
- No specific software installed on clusters: Based on linux **system commands** (bash, ssh, scp, tar, rsync, ...)
- Integrated to function with OAR resource management system
- **Modular Architecture:** Easy for development and research

OAR: A cluster resource management system with high level components

- Simple, robust and scalable
- High level tools (scripting language, relational database engine)
- Designed to be usable as a production resource management system (ex: GRID5000)
- Performance similar with complex resource management systems (LSF,Torque+Maui,SGE,...)

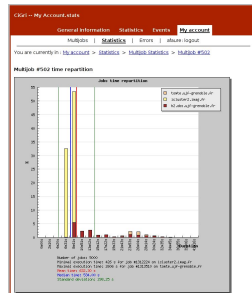
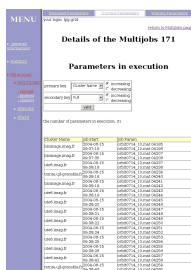
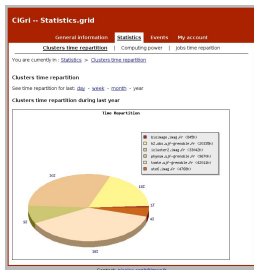
[Ref: Nicolas Capit et Al.(2005)

A batch scheduler with high level components.(CCGrid05)]



CIGRI Grid Functionalities

• Web Portal for grid monitoring (+statistics)

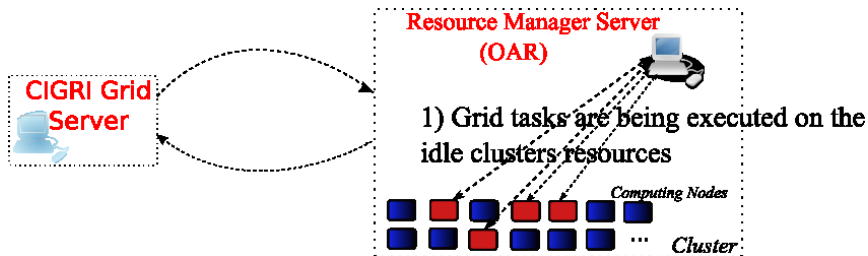


- **Collection of results** of the executed jobs from the clusters on a centralised server (scp)
- Automatic **Data synchronization** of clusters (rsync)
- Support of **diskless** PCs environment
(<http://computemode.imag.fr>)

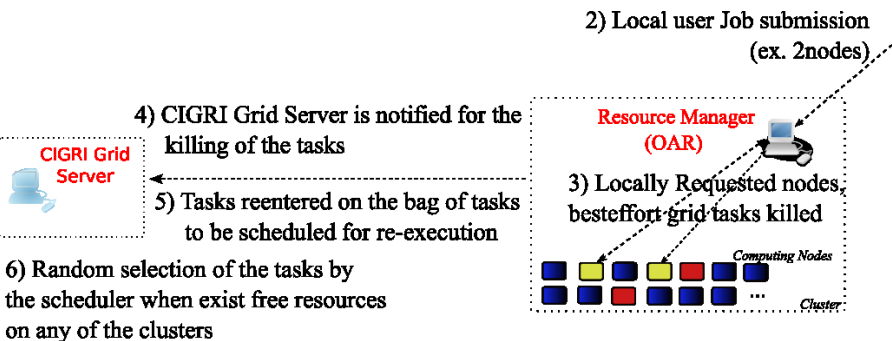
Plan

- 1 CIGRI: lightweight grid computing approach
- 2 **Fault Treatment and Turnaround time Optimisations**
 - CIGRI Fault Treatment
 - Checkpoint/Restart techniques for turnaround time optimisation
- 3 Evaluating CIGRI upon Grid5000 platform
- 4 Conclusions

CIGRI Fault treatment



Fault treatment - No Checkpoints Strategy



- CIGRI guarantees the complete execution of the application

Motivations for optimisation

Why optimise?

- Valuable computation (of hours) can become completely useless just after an interference failure (tasks have to start from the beginning)
- CIGRI default Fault Treatment Strategy cannot guarantee a fast turnaround time.

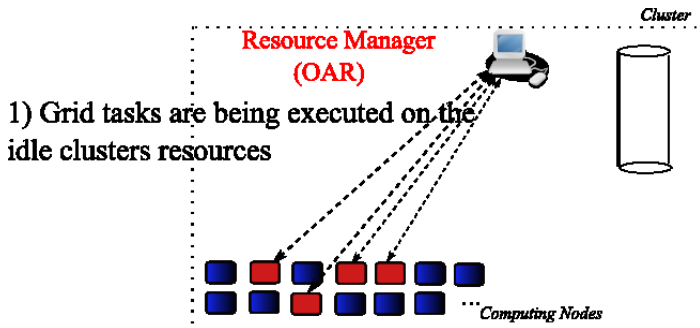
Solution:

- **Checkpoint/Restart** Technique: to cope with the high failure rate of computing nodes.
- Types: Application level/System level (BLCR library)

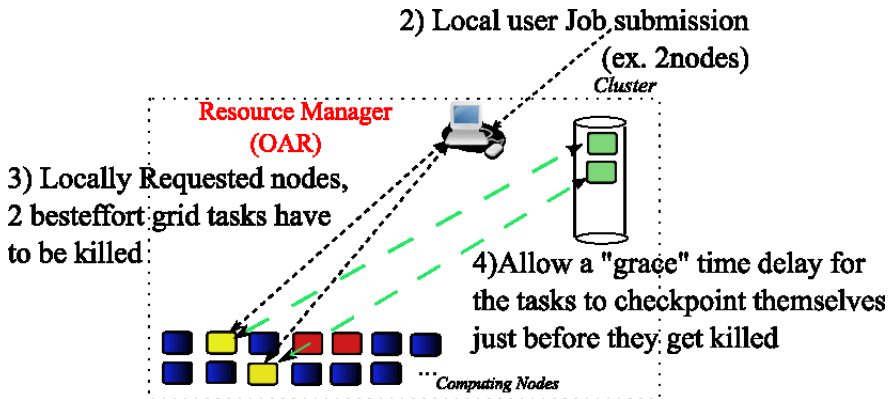
Constraints:

- Application level: Need to change the application code
- System level: BLCR limitations (TCP/UDP Sockets, open file locks, asynchronous I/O not supported)

Triggered checkpoints strategy



Triggered checkpoints strategy



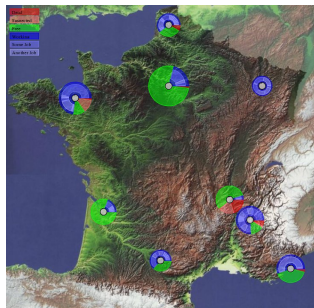
- Used only for interference failures
- Drawback: "grace" time delay

Plan

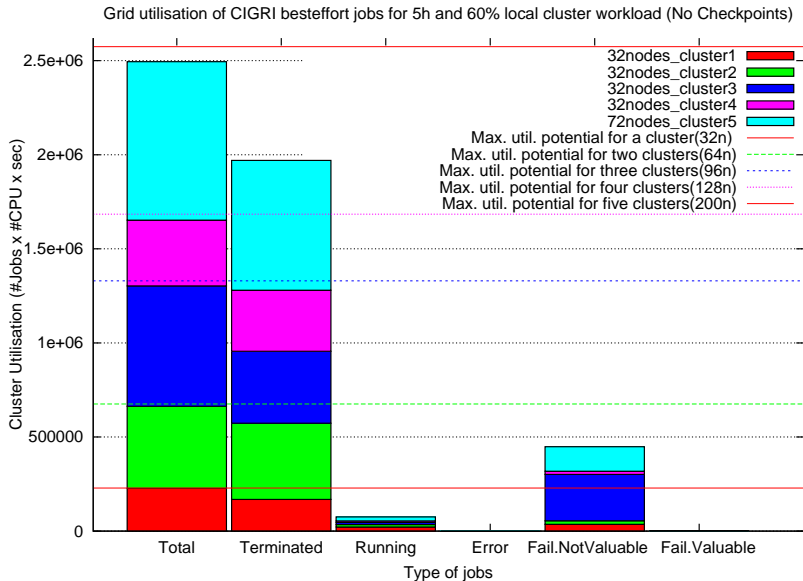
- 1 CIGRI: lightweight grid computing approach
- 2 Fault Treatment and Turnaround time Optimisations
- 3 Evaluating CIGRI upon Grid5000 platform
- 4 Conclusions

Real-life reproducible experiments upon Grid5000

- 1 A CIGRI grid deployed upon Grid5000.
- 2 Local cluster jobs based on obtained grid traces (sleep jobs, no computation just interference).
- 3 One real-life multiparameter application send to CIGRI (real computation).
- 4 Evaluate OAR/CIGRI platform (ex. benchmark the **Fault-tolerance** mechanism and optimisations).
- 5 Collection of results and a posteriori treatment.

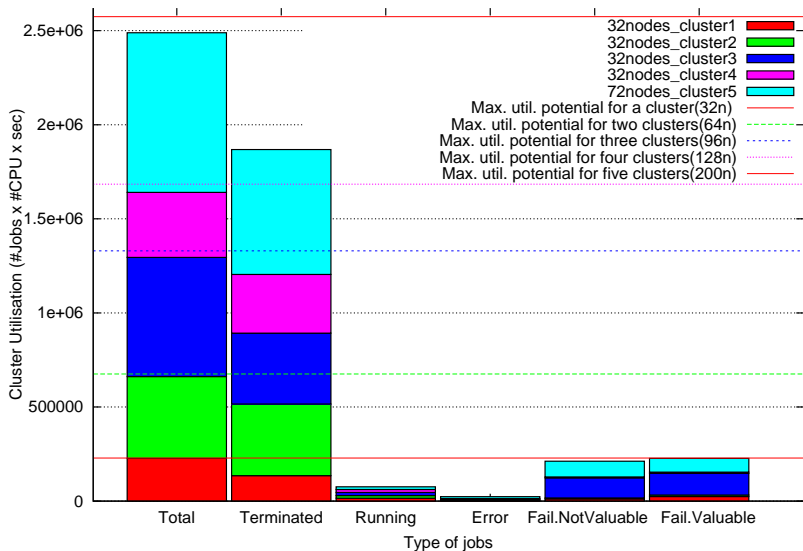


CIGRI grid utilisation for 5 clusters of 200nodes grid, No checkpoints strategy (default)



CIGRI grid utilisation for 5 clusters of 200nodes grid, Triggered checkpoints strategy

Grid utilisation of CIGRI besteffort jobs for 5h and 60% local cluster workload (Triggered Checkpoints)



CIGRI grid experiment Results

- State of grid jobs for 5hours of experimentation of 5 clusters, 200nodes(DUAL OPTERON 2.0GHz, 2GB RAM), 60% local workload

Strategies / Jobs	Total	Terminated	Running	Error	Inter. Failures Not Valuable	Inter. Failure Valuable
Triggered checkpoints	1377	762	74	103	226	212
No checkpoints	1420	739	74	8	581	0

[Ref: Yiannis Georgiou et Al.(2007)

Evaluations of the Lightweight Grid CIGRI upon the Grid5000 Platform (eScience2007)]

Plan

- 1 CIGRI: lightweight grid computing approach
- 2 Fault Treatment and Turnaround time Optimisations
- 3 Evaluating CIGRI upon Grid5000 platform
- 4 Conclusions**

Conclusions

CIGRI Principal contributions

- A **Lightweight** grid computing approach
- CIGRI Grid jobs: **transparent** execution upon the interconnected clusters.
- **Best-effort** jobs and the utilisation of unused resources
- Checkpoint/Restart for faster turnaround time

CIGRI-EGEE Interactions

- 1 Transparency of Best-effort grid jobs on EGEE...?? Could it be valuable?
- 2 Current project: Cohabitation of CIGRI RhoneAlpes Grid with gLite EGEE Grid
 - Usage: Local cluster jobs, Low priority besteffort grid jobs(CIGRI) and remote grid jobs(EGEE)
 - OAR Resource Management System can cohabitate with any RMS system (SGE,LSF,PBS,..) connected with gLite



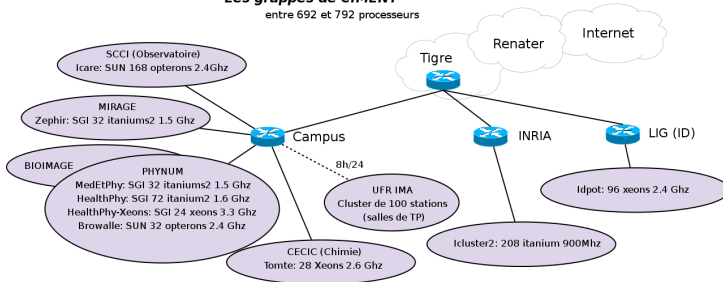
CiGri



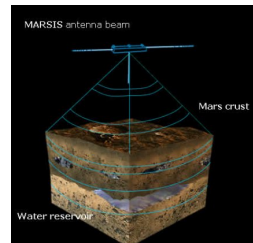
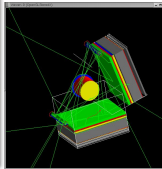
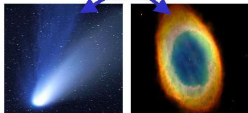
CIMENT: an OAR/CIGRI grid research platform

Les grappes de CIMENT

entre 692 et 792 processeurs



H₂O



Links

CiGri

<http://cigri.imag.fr>

<http://cigri.ujf-grenoble.fr>



<http://oar.imag.fr>

CIGRI vs gLite: Principal Differences

CiGri

- Ideal for Regional or Metropolitan Area Grids
- Basic Security
- Simple (ssh, bash, rsync)
- Support of only BoT applications
- Best-effort jobs and the utilisation of unused cluster resources



- World Area Grid
- Strong security measures
- Standards (GRAM, GridFTP)
- Support of BoT + MPI + DAG Dependent applications
- Dedicated cluster resources
- Virtual Organizations