### Cluster Computing Resource and Job Management for HPC

SC-CAMP

16/08/2010







3 🕨 🖌 3

(SC-CAMP)

**Cluster Computing** 

16/08/2010 1 / 50

### Summary

- Introduction Cluster Computing
- 2 About Resource and Job Management Systems
- 3 OAR a highly configurable RJMS
  - OAR Concepts and architecture
- OAR Scheduling
- OAR Interfaces
  - 7 Research Issues and Conclusions

★ ∃ ►

Introduction Cluster Computing

#### Top500 larger computing infrastructures

Architecture / Systems June 2010



http ://www.top500.org/

(SC-CAMP)

イロト イヨト イヨト イヨト

#### Definition Cluster

A computer cluster for High Performance Computing is a group of linked computers, working together closely connected to each other through fast local area networks. A cluster is also called *parallel supercomputer*.



- 金属 トーイ

#### Process

Processes are runing on CPUs.

- A process is a program that is running into memory
- On UNIX, several processes may be running on one or several processors (multitask). They are hierarchical and belong to a particular user.
- Each UNIX process has a unique number on a given system. It is called the PID.
- A "thread" (also known as a "light process") is a program that shares a certain amount of memory with another thread belonging to the same father process (under Linux, threads are viewed with "ps -eLf")



→ Ξ →

#### Jobs

Processes can be groupped into jobs.

- A *job* may be a single process, a group of processes or even a batch.
- In our context (HPC clusters), a job is a set of processes that have been automatically launched by a scheduler by the way of a user's submitted script.
- A job may result in N instances of the same program on N nodes or processors of a cluster.



#### Nodes

Jobs are running on nodes.

• A *node* is a computer having a set of p CPU, an amount of memory, one or several network interfaces and that may have a storage unit (local disk).



- (三) - (-)

#### Computing network

Several nodes are connected to a computing network, generally low latency network (Myrinet, Infiniband, Numalink,...)



#### Types of jobs

Jobs maybe "parallel" or "sequential". A parallel job runs on several nodes, using the computing network to communicate.

- Sequential Job is a unique process that runs on a unique processor of a unique host.
- Parallel Job Several processes or threads that can communicate via a specific library (MPI, OpenMP, threads,...). Some of them are 'shared memory' parallel jobs (they run on a unique multiprocessor host) and some are 'distributed memory' parallel jobs (they may run on several hosts that communicate via a low latency high speed network) Warning : a 'fake' parallel job may hide several sequential jobs (several independent processes)



#### Types of jobs

The Jobs can come at the form of Batch or Interactive type.

- A *batch* job is a script. A shell script is a given list of shell commands to be executed in a given order
  - Todays shells are so sophisticated that you can create scripts that are real programs with some variables, controle structures, loops,
  - We sometimes also call script a program that has been written in an interpreted language (like perl, php, or ruby)
- An *interactive* job is usually an allocation upon one or more nodes where the user gets a shell on one of the cluster nodes.

A (10) A (10) A (10)

#### Resource and Job Management System

The *Resource and Job Management System* (RJMS) or Batch Scheduler is a particular software of the cluster that is responsible to distribute computing power to user jobs within a parallel computing infrastructure.



(SC-CAMP)

#### Concepts Resource and Job Management System

The goal of a Resource and Job Management System is to satisfy users demands for computation and achieve a good performance in overall system's utilization by efficiently assigning jobs to resources.



12 / 50

#### Concepts Resource and Job Management System

This assignement involves three principal abstraction layers :

- the declaration of a job where the demand of resources and job characteristics take place,
- the scheduling of the jobs upon the resources
- and the launching and placement of job instances upon the computation resources along with the job's control of execution.

In this sense, the work of a RJMS can be decomposed into three main subsystems : Job Management, Scheduling and Resource Management.

RJMS subsystems	Principal Concepts	Advanced Features
Resource Management		
	<ul> <li>Resource Treatment (hierarchy, partitions,)</li> </ul>	- High Availability
	-Job Launcing, Propagation, Execution control	<ul> <li>Energy Efficiency</li> </ul>
	<ul> <li>Task Placement (topology, binding,)</li> </ul>	<ul> <li>Topology aware placement</li> </ul>
Job Management		
	<ul> <li>Job declaration (types, characteristics,)</li> </ul>	<ul> <li>Authentication (limitations, security,)</li> </ul>
	<ul> <li>Job Control (signaling, reprioritizing,)</li> </ul>	<ul> <li>QOS (checkpoint, suspend, accounting,)</li> </ul>
	<ul> <li>Monitoring (reporting, visualization,)</li> </ul>	<ul> <li>Interfacing (MPI libraries, debuggers, APIs,)</li> </ul>
Scheduling		
	-Scheduling Algorithms (builtin, external,)	- Advanced Reservation
	-Queues Management (priorities,multiple,)	- Application Licenses

TAB.: General Principles of a Resource and Job Management System

(日) (同) (日) (日) (日)

Scheduling Policy	Description
<u>FIFO</u>	jobs are treated with the order they arrive.
<u>Backfill</u>	fill up empty wholes in the scheduling tables without modifying the order or the execution of previous submitted jobs.
Gang Scheduling	multiple jobs may be allocated to the same resources and are alternately suspended/resumed letting only one of them at a time have dedicated use of those resources, for a predefined duration.
TimeSharing	multiple jobs may be allocated to the same resources allowing the sharing of computational resources. The sharing is managed by the scheduler of the operating system
<u>Fairshare</u>	take into account past executed jobs of each user and give priorities to users that have been less using the cluster.
Preemption	suspending one or more "low-priority" jobs to let a "high-priority" job run uninterrupted until it completes

#### $\operatorname{TAB.:}$ Common Scheduling policies for RJMS

æ

・ロト ・聞ト ・ヨト ・ヨト

#### Opensource and Commercial RJMS

- Context : High Performance Computing
- A lot of them : Condor , Sun Grid Engine (SGE), MAUI/Torque, Slurm, OAR, Catalina, *LSF* , Lava, *PBS Pro*, *Moab*, Loadleveler, CCS...
- http://en.wikipedia.org/wiki/Job\_scheduler

イロト イヨト イヨト

#### **RJMS** General organization

- A central host
- Client programs (command line) to interract with users
- A lot of configuration parameters



3

(日) (周) (三) (三)

## RJMS Features (1/2)

#### non-exhaustive list

- Interactive tasks (submission)(shell) / Batch
- Sequential tasks
- Walltime limit. (very important for scheduling !)
- Exclusive / non-exclusive access to resources
- Resources matching
- Scripts Epilogue/Prologue (before and after tasks)
- Monitoring of the tasks (resources consumption)
- Job dependencies (workflow)
- Logging and accounting
- Suspend/resume

イロト 不得下 イヨト イヨト

## RJMS Features (2/2)

#### non-exhaustive list

- Array jobs
- First-Fit (Conservative Backfilling,)
- Fairsharing
- ...

(日) (周) (三) (三)

#### Objectives

OAR has been designed with versatililty and customization in mind.

- Following technological evolutions (machines and infrastructures more and more complicated)
- Initial context : CIMENT and Grid5000
- Different contexts adaptation (cluster, cluster-on-demand, virtual cluster, multi-cluster, lightweight grid, experimental platform as Grid'5000, *big cluster*, special needs).



(日) (同) (三) (三)

#### CIMENT local HPC center

- Université Joseph Fourier (Grenoble)
- A dozen of heterogeneous supercomputers, more than 2000 cores today
- special feature : a lightweight grid (CiGri) making profit of unused cpu cycles through best-effort jobs (a OAR feature)
- Great collaborative work production/research

#### Grid 5000

- Grid for computing experiments
- 9 sites in France + 2 sites abroad
- Interconnection gigabit Renater
- More than 5000 cores today

- ∢ ∃ ▶

### OAR special features

#### non-exhaustive list

- Classical features +
- Advance Reservation
- Hierarchical expressions into requests
- Different types of resources (ex licence, storage capacity, network capacity...)
- Besteffort tasks (zero priority task, highly used by CiGri)
- **Multiple task types** (besteffort, deploy, timesharing, idempotent, power, cosystem ...) (customizable)
- Moldable tasks
- Energy saving

(日) (周) (三) (三)

### OAR : design concepts

High level components use

- **Relational database** (MySql/PostgreSQL) as the kernel to store and exchange data
  - resources and tasks data
  - internal state of the system

• Script languages (Perl, Ruby) for the execution engine and modules

- Well suited for system parts of the code
- High level strcutures (lists, hash tables, sort...)
- Short developpement cycles
- Other components : SSH, CPUSETS, Taktuk
  - SSH, CPUSET (isolation, cleaning)
  - Taktuk adaptative parallel launcher

(人間) トイヨト イヨト

#### OAR : general organisation



16/08/2010 25 / 50

#### Job life cycle



10	$\sim$	~	A 1		0
5	C - '	с.	Δ.	M	P
, Ξ	~	~			• •

16/08/2010 26 / 50

3

<ロ> (日) (日) (日) (日) (日)

#### Task status diagram



(SC-CAMP)

**Cluster Computing** 

16/08/2010 27 / 50

### Submission examples : OAR

Interactive task : 1

oarsub -l nodes=4 -l

Batch submission (with walltime and choice of the queue) :

 oarsub -q default -l walltime=2 :00,nodes=10 /home/toto/script

Advance reservation :

```
• oarsub -r "2008-04-27 11 :00" -l nodes=12
```

Connection to a reservation (using task's id) :

```
• oarsub -C 154
```

<sup>1</sup>Note : Each submission command returns an id number.  $\bullet \triangleleft \bullet \bullet \bullet \models \bullet$ 

#### **Hierarchical Resources**



### Scheduling

- Scheduling is the process that is responsible to assign jobs according to the users needs and predefined rules and policies, upon available computational resources that match with the demands.
- A typical scheduler functions in cooperation with queues which are elements defined to organize jobs according to similar characteristics (for example priorities).
- Numerous **parameters** are used to guide and scope the allocations and priorities.

<sup>&</sup>lt;sup>2</sup>Note : scheduling is computed again at each major change in the state of a task. O ( (SC-CAMP) Cluster Computing 16/08/2010 30 / 50

### Scheduling organization into OAR

#### Tasks are put into queues

- each queue as a priority level
- each queue obey to a scheduling policy



#### Resource matching

#### A preliminary step to scheduling

- Resources filtering
- Resources sorting
- Allows the user to have special needs
- memory, architecture, special hosts, OS, workload...

### FIFO : Fisrt-In First-Out



1	0	$\sim$	~			5	
L	5	('	L.,	4 I	VΙ	$\mathbf{P}$	
۰.						• •	

2

<ロ> (日) (日) (日) (日) (日)

## First-Fit (Backfilling)

#### Holes are filled if the previous tasks order is not changed.



16/08/2010

|田・ (日) (日

34 / 50

### FairSharing

The order of tasks is computed depending on what has already been consumed (low time-consuming users are prioritized) A time-window and weighting parameters are defined.



(SC-CAMP)

16/08/2010 35 / 50

### TimeSharing



1 1	~ .	~ .	~ .			<b>n</b> )
( )	50	(	_/	41	VI.	Ρ.

16/08/2010 36 / 50

- 2

<ロ> (日) (日) (日) (日) (日)

#### Advance Reservation

- Very useful for demos, planification, multi-site tasks or grid tasks...
- But :
  - Restrictive for the scheduling
  - Resources are rarely used along the whole duration of the reservation (waste!)



#### OAR : Monika

(SC-CAMP)

#### **OAR** Cluster nodes

default summary			
	Free	Busy	Total
network_address	4	15	32
resource_id	32	120	256

#### **Reservations:**

Reservations for property iru=0:

	1000000		1000 00	100070	100070	100070					
<u>r1i0n0</u>	182270	182270	182270	182270	182270	182270	182270	182270			
rliOnl	182270	182270	182270	182270	182270	182270	182270	182270			
r1i0n2	Free	Free	Free	Free	Free	Free	Free	Free			
<u>r1i0n3</u>	Free	Free	Free	Free	Free	Free	Free	Free			
<u>r1i0n4</u>	182267	182267	182267	182267	182267	182267	182267	182267			
<u>rli0n5</u>	Free	Free	Free	Free	Free	Free	Free	Free			
<u>r1i0n6</u>	<b>182271</b>	182271	182271	182271	182271	182271	182271	182271			
<u>r1i0n7</u>	182271	182271	182271	182271	182271	182271	182271	182271			
r1i0n8	182271	182271	182271	182271	182271	182271	182271	182271			
<u>r1i0n9</u>	StandBy	StandBy	StandBy	StandBy	StandBy	StandBy	StandBy	StandBy			
<u>r1i0n10</u>	StandBy	StandBy	StandBy	StandBy	StandBy	StandBy	StandBy	StandBy			
<u>r1i0n11</u>	182294	182294	182294	182294	182294	182294	182294	182294			
r1i0n12	182282	182282	182282	182282	182282	182282	182282	182282	三国 とう 国 と	臣	9
			Clus	ter Com	puting				16/08/20	10	38

#### OAR : Gantt diagramm

🗣 + 📦 + 🎒 🙆 😤	http://ita101.imag.fr/cgi-bin/DrawOARGant	t.pl
---------------	---	------

💌 🔘 ок [ \_ @ × Origin 2005 \_ Oct \_ 17 \_ 00:00 \_ Range 3 days \_ I BestEffort \_ Draw \_ Default

(SC-CAMP)

**Cluster Computing** 

16/08/2010 39 / 50

#### **Research Issues**

Topological Constraints for improvement of application performance.Energy consumption mechanisms in HPC clusters

### Topological constraints

#### Hardware evolution

- switch/node/cpu/core : Hierarchical architecture
- NUMA host/ BlueGene host : 2D grid, 3D or hybrid •



#### 4 x 2 CPUs / 4 COREs

-	<u> </u>	-			>	
5	( _		Δn	<i>.</i> /I	$\mathbf{P}$	a
ິ	<u> </u>	$\sim$		•••		

3

• • = • • = •

### Hierarchical topological constraints

Problem with parallel applications that are network-bandwidth sensitive.



4 x 2 CPUs / 4 COREs

$(SC_C \Delta MP)$

3

- 4 同 6 4 日 6 4 日 6

### Hierarchical topological constraints

- Hierarchy in the requests : oarsub -I switch=1/nodes=2/cpu=2/core=2 ./my-app = 1x2x2x2 = 8 cores
- Linux cpusets management (be aware of the cpu affinity inside a cpuset)

(日) (周) (三) (三)

#### Energy saving

- We can shutdown nodes when not used (on-demand wake up)
- Timetable priorities
- Work done during Google Summer Of Code 2008 (Gsoc'08)
  - Parametrical job type : **powersaving** + options (cpufreq, disk shutdown, video ..., special policy)
  - $\bullet~\mbox{Ex Job BestEffort} \rightarrow \mbox{Iowest CPU frequency.}$



### Conclusion

- Management of Resources and Jobs for high performance computing in modern architectures has become a complex procedure with interesting research issues.
- Opensource and commercial RJMS have been evolving to provide efficient exploitation of the infrastructures along with quality of services to the users.
- OAR versatile and customizable distributed RJMS.

# Questions?



http://oar.imag.fr/

		<u> </u>	-			
1	5	( _		Δ	NΛ	P
١.	-	<u> </u>	<u> </u>			

3

・ロン ・四 ・ ・ ヨン ・ ヨン

#### Liens



Condor http://www.cs.wisc.edu/condor/

- Sun Grid Engine (SGE) http://gridengine.sunsource.net
- TORQUE/MAUI http://www.clusterresources.com/



SLURM

www.llnl.gov/linux/slurm/



http://www.platform.com



OAR http://oar.imag.fr

• • = • • = •



#### Energy consumption for trace file execution of 50.32% system utilization

(SC-CAMP)

#### **Cluster Computing**

16/08/2010 48 / 50



#### Energy consumption of trace file execution with 89.62% of system utilization



(SC-CAMP)

æ

<ロ> (日) (日) (日) (日) (日)