

Tasks Mapping with Quality of Service for Coarse Grain Parallel Applications

Patricia Pascal, Samuel Richard, Bernard Miegemolle, and Thierry Monteil

LAAS-CNRS, 7 Avenue du Colonel Roche 31077 Toulouse France
{ppascal,srichard,bmiegemo,monteil}@laas.fr

Abstract. Clusters and computational grids are opened environments on which a great number of different users can submit computational requests. Some privileged users may have strong Quality of Service requirements whereas others may be less demanding. Common mapping algorithms are not well suited to guarantee a defined quality of service, they propose at best priority systems in order to favour some applications without any guaranty. We propose a new mapping algorithm, dealing with the notion of quality of service for scheduling applications over clusters and grids over different classes of service.

This algorithm uses information on the application to map, all the unfinished applications previously mapped, the state of the execution support, and the processor access model (round robin model) to suggest a mapping which guarantees all the expressed constraints. The mapping decision is taken on-line based on the release date of all applications and the memory space used. To finish, the validation of the algorithm is performed with real log files entries simulated with Simgrid.

Keywords: scheduling, quality of service, resource manager, grid, clusters

1 Introduction

In distributed environments, resource management is very important in order to take advantage of multiple hosts and to optimize resource use. The scheduling policy commonly used on distributed systems is best effort with priorities. Different queues are created: short jobs, long jobs, high parallel jobs, etc with FIFO or more elaborated policies. This system of queues can be used to allow the differentiation of users by assigning a priority to each queue but does not guarantee any quality of service. This limitation is due to historical reasons because batch schedulers have been created for parallel computers which are generally used by few users. Clusters are more opened and also complicated environments. Due to their low cost, they can be accessed by a lot of different users that have different needs and expectations; they are also connected with network on which different policies of quality of service can be used. For this reason, batch schedulers are not well suited to ensure different qualities of service to many users using the same execution environment. In this article, a scheduling algorithm, used in distributed systems like clusters or aggregation of clusters (grid), and

which implements different classes of service, is presented. This algorithm is implemented in a tool called AROMA (scAlable ResOurces Manager and wAtcher) [4]. AROMA integrates a resource management system, an application launcher, a scheduler, a statistic module and an accounting system.

In the first section, a state of the art is presented; then the context of the study, the notion of quality of service and AROMA are detailed. After that, the optimization problem that has been solved is explained. To conclude, first results validating the proposed algorithm are given. The originality of this work is to mix different processes from different classes of service on the same processor at the same time.

2 Related Work

Batch and dynamic schedulings are difficult problems to solve because resources needed by an application may not be known. Resources are heterogeneous and their availability is not completely known. Moreover, the mapping algorithm must run quickly, therefore heuristics with good properties are used. Henri Casanova in [5] studies deadline scheduling on computational grid. His goal is to minimize the overall occurrences of deadline misses as well as their magnitude.

Rajkumar Buyya in [6] proposes a deadline and budget constrained cost-time optimization algorithm for scheduling on grids. The algorithm is called DBC (Deadline and Budget Constrained).

Mechanisms have been created to improve the mapping. The first one concerns resource reservation. Different types of resource reservation algorithms are studied in [1]. They evaluate the performance with or without preemption. The reservation insures that all the resources necessary to run the applications will be free. A second mechanism is the gang-scheduling. It creates time slices, that is to say parts of time that are allocated to the parallel and sequential applications [2]. With this solution, all applications progress simultaneously. Nevertheless, it could create a problem of overload and memory saturation. Finally, the backfilling [3] allows the insertion of jobs into scheduler queue. The insertion is possible if it does not perturb the other jobs. It is a way to remove the holes in resources utilization.

This article proposes a way to mix jobs requirements with different qualities of service (deadline, immediate execution, dedicated resources).

3 Mapping Algorithm

As the context of the study is ASP (Application Service Provider), the hypothesis that all the applications consuming resources are known is made : that is to say, hosts are considered dedicated to computation. All the jobs are submitted through AROMA and system tasks influence is neglected. AROMA daemons are also able to monitor running applications; this information is used to refresh estimated completion date of running jobs. The second hypothesis made

is that applications are regular coarse grain parallel applications for which the time spent in communication is small and the execution time can be roughly predicted. The problem is to find the mapping of a new application knowing all the previously mapped applications that are still consuming resources in the system. The proposed mapping has to guarantee that the quality of service is respected for all applications (running and currently scheduled applications).

3.1 Quality of Service and Mapping Problem

Applications are grouped into four application classes (in order of importance):

- **Deadline applications (class 1):** this class of service guarantees that the execution will end before the deadline. Execution can be immediate or deferred.
- **High priority applications (class 2):** this class of service guarantees that the execution will be immediate.
- **Applications with dedicated resources (class 3):** this class of service guarantees that each application will be the only one to use resources during its execution. Execution can be immediate or deferred.
- **Applications without constraint (class 4):** this class of service corresponds to applications which will be executed as soon as possible with available resources. This class is also named “Best Effort”.

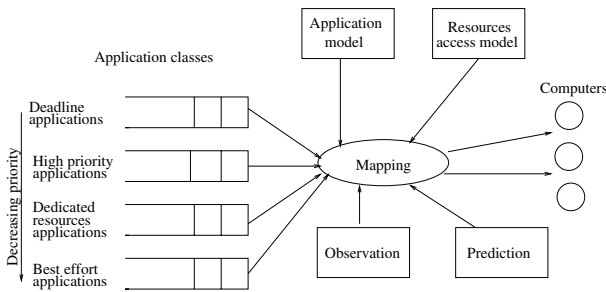


Fig. 1. The mapping problem

The mapping problem inputs are (figure 1):

- **Application model:** Some information describing the application needs and requirements has to be supplied to the scheduler in order to take a good mapping decision. Some elements are inputs of the algorithm while others express constraints for the mapping. Inputs are an estimation of the cpu time required by each task, the number of tasks and the size of exchanged data between the different tasks of the application. Those values can be given by the user or retrieve

from a database containing information on previous runs for the same type of application. The application model can express additional constraints like the fact that all the tasks must begin at the same date or temporal relations between the tasks (classical description in graph theory).

Software or specific hardware requirements add constraints on the execution hosts. Class 1 induces a constraint on the ending date of the application, class 2 induces a constraint on the starting date of the application and class 3 induces a constraint on the execution hosts.

All the constraints are verified by the algorithm.

– **Resources access model:**

According to the resources access policy, equations are deduced and used to predict the utilization time of the resources and the end of execution. Models developed in this article are deterministic, nevertheless models which take care of random perturbation (arrival of uncontrolled jobs, for example a direct login on the host) have been developed in [12]. The difficulty is to mix different applications from different classes of service on the same host while respecting all the constraints: several tasks may share the same host during the same period of time.

– **Observations:**

The processors and network load (percentage of processors utilization, bandwidth used), idle memory space and number of processes are monitored. They are used to refresh information used for mapping.

– **Mapping:**

Different queues exist, each corresponding to a priority level (figure 1). When several applications have to be mapped, the jobs of the highest priority queue will be treated first. If no mapping respecting the constraints is found, three cases are studied:

- the algorithm try to move a job with a weaker priority : it looks for an application in a lower priority queue which has been planned to be run in the future, then it removes it and try to map it again after the current mapping.
- if the previous case is impossible, the algorithm can stop a running application and try to find a new mapping for this application. In this case, a new constraint is created to express that this application must go on later on the same host. There is no migration.
- if the two previous cases are impossible, the mapping request is rejected.

3.2 Mathematical Expression of the Problem

The Variables

- t_0 : initial date of the mapping research.
- $t_b(a, p, m)$: execution starting date of the task p of the application a on the host m .
- $t_{fn}(a, p, m)$: end of execution of the task p of the application a on the host m when the network is neglected.

- $t_f(a, p, m)$: end of execution of the task p of the application a on the host m when an estimation of time spent on communication is done.
- $t_c(a, p, m)$: processor time requested by the task p of the application a on the host m . Coarse grain applications are considered, so “small” communication time is taken into account and the synchronization time is neglected.
- $t_c^k(a, p, m)$: remaining processor time for the task p of the application a after the event number k on the host m (events are a beginning or the end of a task).
- $D_r(a, p)$ estimation of the size of data sent and received by the task p of the application a .
- $C(m)$: coefficient to take care of heterogeneous processors.
- $t_c(a, p, m) = t_c(a, p) * C(m)$: equivalence of processor time requested by the task p of the application a for the host m .
- $B(i, j)$: estimation of bandwidth of the network between host i and host j .
- M : number of hosts. A : number of applications to map.
- N^a : number of tasks of application a .
- $X_m(t)$: number of processes on the host m at time t .
- t_m^k : a beginning event or an ending event of a process on the host m . It is the date of the event number k .
- $P(a)$: set of possible mapping for application a .
- $t_f(m, s)$: the end of all the tasks mapped on the machine m after the mapping s of application a with $s \in P(a)$.
- $M(a, p)$: the machine on which the task p of the application a is executed.
- $U(m)$: number of processors available on host m .
- $M_u(m)$: total memory used on host m . $M_t(m)$: total memory on host m .
- M_f : constant to modify the weight of the memory criteria.

The Optimization Criteria: The mapping problem is an optimization problem, criterion has to be chosen. Several criteria are well known ([7][8][9][10]) : makespan (minimizing the termination date of an application), sum-flow (minimizing the quantity of resources used), max-stretch (it expresses that a task has been slowed compared to what its execution would have been on an idle server). The objective here is to optimize the use of the providers resources and to guarantee the level of quality of service required. So it has been chosen to liberate all the resources as soon as possible. Applications already mapped may be influenced by the mapping found. By consequence, all the applications (the currently mapped and the previously mapped) must finish as soon as possible.

The date of the end of resources utilization is optimized. Moreover a second criteria consists in moderating with memory space used: hosts which have the most free memory space are privileged first. The choice has been made to introduce memory in criteria because the exact amount of memory requested by an application is often unknown by users. The problem is multi-criteria by using a linear combination of different criteria (1). The first part of the addition ($t_f(m, s)$) refers to the release date of the machine. The second part of the addition $t_f(m, s) * M_u(m) * M_f/M_t(m)$ penalizes machines which have less free

memory. $M_u(m)/M_t(m)$ gives an idea of memory utilization on this host. The coefficient M_f influences the weight of this part of criterion. The multiplication with $t_f(m, s)$ puts the value into the same order of value as the first part of the criterion.

$$\min_{s \in P(a)} (\max_m (t_f(m, s) + t_f(m, s) * M_u(m) * M_f / M_t(m))) \quad (1)$$

The Mapping Algorithm: A list algorithm for applications, tasks and machines is used to reduce the combinatorial. The mapping of an already studied task of an application a is revised only if it is impossible to find a mapping for this application. Moreover, in order to reduce the time used by the algorithm, for each host, release date is saved and hosts are ordered to study which of them will give the best mapping first.

The quality of service is already respected, all the constraints induced are verified by the algorithm. If it is impossible to find a mapping corresponding to the demand, the request is refused.

The computation of $t_b(a, p, m)$ and $t_f(a, p, m)$ will now be explained. The equations are found considering that the processor access follows a round robin policy.

The Starting Date of a Task: $t_b(a, p, m)$ is computed with an iterative algorithm (at the beginning it is t_0). If, at this date, no mapping can be found, another date is searched.

$$\begin{aligned} t_b^0(a, p, m) &= t_0 \\ t_b^{k+1}(a, p, m) &= \min t_f(i, j, m) \end{aligned}$$

with $i \in [1, a - 1]$, $j \in [1, N^i]$ and $t_f(i, j, m) > t_b^k(a, p, m)$ (only the tasks that can end after the last t_b studied, are considered). The idea is to search a new starting date when the system is less loaded: when a job finishes.

The End of a Task: $t_f(a, p, m)$ is computed with an iterative algorithm. Each date is studied when there is a creation or a termination of a job. t_m^0 corresponds to the arrival of the first process on the host.

for all tasks (a goes from 1 to A and p from 1 to N^a , on m)

if $X_m(t_m^k) > U(m)$ (is there more processes than processors ?)

$$t_m^{k+1} = \min_{a,p} (t_b(a, p, m), t_c^k(a,p,m) * X_m(t_m^k) / U(m) + t_m^k) \quad (\text{next event corresponds to a creation or a death of a process})$$

else

$t_m^{k+1} = \min_{a,p} (t_b(a, p, m), t_c^k(a,p,m) + t_m^k)$ (next event corresponds to a creation or a death of a process)

if $t_b(a, p, m) > t_m^k$ and if $X_m(t_m^k) > U(m)$

$t_c^{k+1}(a, p, m) = t_c^k(a, p, m) - U(m) * (t_m^{k+1} - t_m^k) / X_m(t_m^k)$ (estimation of the new requested time of processor for this task after this short execution on processor: it is the time sharing policy)

$$\begin{aligned} \text{if } t_b(a, p, m) = t_m^i \\ t_c^i(a, p, m) = t_c(a, p, m) \\ \text{(it is case of an insertion of a new process in the recurrence)} \end{aligned}$$

Estimation of the finished date of process without the network:

$$\text{if } t_c^{k+1}(a, p, m) = 0, t_{fwn}(a, p, m) = t_m^{k+1}$$

Time spent in communication are put inside the estimation of the end of process to advantage location of tasks on the same cluster or on the same site because the bandwidth will be better. The unmapped tasks are ignored in the estimation of the worst bandwidth used for the application a .

$$t_f(a, p, m) = t_{fwn}(a, p, m) + D_r(a, p) / \min_{i \in [1, N^a]} B(M(a, i), m)$$

The iterations continue until all the tasks of the host finish. In fact, the mapping algorithm quickly simulates the execution of jobs and can mix the different classes of service.

4 Validation

To validate the algorithm, Simgrid ([11]) simulator has been used. Real jobs submission log files have been used to estimate the behavior of the algorithm with Simgrid. Nevertheless, it is difficult to compare the algorithm to others because algorithms found do not define classes of service and do not execute processes of the different classes at the same time on the same processors.

4.1 Comparison with NQS

Feitelson logs (real logs) have been used. These logs give : the submission date of jobs, the required cpu time, the number of tasks. The log file *Lsdsc-sp2.swf* [13] is used. This Job Trace Repository is brought by the HPC Systems group of the San Diego Supercomputer Center (SDSC), which is the leading-edge site of the National Partnership for Advanced Computational Infrastructure (NPACI) [14], [15]. The real system has 128 nodes and is scheduled with NQS [16]. Jobs submissions are reproduced, the mapping research is done with the algorithms on 128 nodes and their execution is simulated with Simgrid. The mean waiting times given by the logs are compared to the mean waiting times of the algorithms with quality of service. So, 10000 and 35000 jobs are simulated. Each job was synchronized, this means that all the tasks of the same parallel application must begin at the same date and belong to the dedicated resources class. This class seems to be the nearest from NQS policy. The simulation of 10000 jobs is equivalent to an activity on the supercomputer during 112 days. The simulation of 35000 jobs is equivalent to an activity on the supercomputer during 249 days. There is no information about the communications, so they are neglected. The results show that, with the proposed algorithm, a better waiting time than NQS is obtained. In fact, with the proposed algorithm, an application can be mapped between two others because the processor time required is known and the prediction of the end of each task mapped can be done. This reduces significantly

the waiting times. NQS sorts the applications into queues (based on required cpu time) and mixes the applications when a mapping is researched. This can increase the waiting time of short jobs which can be slowed by long ones. This problem is avoided with the proposed algorithm because the applications are considered in the order of submission and classes. The second reason is that the proposed algorithm uses more precise values for the requested time of cpu than NQS; by consequence, the mapping is more accurate. The comparison would be more fair if NQS had as much queues than the different times of processor requested. The results are presented in the table 1. The times are given in seconds (s).

Table 1. Comparison of the waiting time between NQS and the proposed algorithm

	10000 events	35000 events
mean waiting time with NQS(s)	10796	8979
mean waiting time with the algorithm(s)	4008	4202

4.2 Influence of Quality of Service

The same log file as previously has been used. In this log file, it is specified that there are four queues (low, normal, high, express), and for each job the queue of submission is known. So this information is used to make an arbitrary correspondence with the proposed classes of applications. The logs are used only to have an approximation of a realistic incoming rate of applications and a good sample of applications requested cpu time. So queue *low* corresponds to best effort class (class 4), queue *normal* corresponds to deadline class (class 1), queue *high* corresponds to dedicated resources class (class 3) and queue *express* corresponds to high priority class (class 2). For deadline applications (class 1), the deadline is : $submission\ date + 5 * cpu\ time\ required$. For high priority applications (class 2), the starting date that must be respected is : $submission\ date + 5\ seconds$. Four cases have been simulated:

- case 1 (10000 events using the four classes) : the waiting time for each class, the global waiting time, the mapping time for each class and the global mapping time are computed.
- case 2 (10000 events in the best effort class) : the global waiting time and the global mapping time are computed.
- case 3 (35000 events using the four classes) : the waiting time for each class, the global waiting time, the mapping time for each class and the global mapping time are given.
- case 4 (35000 events in the best effort class) : the global waiting time and the global mapping time are evaluated.

Doing so, case 1 can be compared with case 2 and then case 3 can be compared with case 4 to see the influence of the quality of service on the mapping

Table 2. Influence of the quality of service on the performance of the algorithm

	Case 1	Case 2	Case 3	Case 4
waiting time class 1 (s)	30.51		30.7	
waiting time class 2 (s)	0.0036		0.003	
waiting time class 3 (s)	4845		5182	
waiting time class 4 (s)	4.55		30.5	
global waiting time (s)	710.5	64.68	396	70.3
mapping time class 1 (ms)	20.22		21.1	
mapping time class 2 (ms)	3.6		3	
mapping time class 3 (ms)	269.3		334.4	
mapping time class 4 (ms)	7.2		28.6	
global mapping time (ms)	50.3	13	49.1	76.4

performances. The results are presented in the table 2. The waiting times are given in seconds (s) and the mapping times in milliseconds (ms).

The introduction of the quality of service increases the mapping times and waiting times but the algorithm still has good performances (mapping time is inferior to 80 ms). The mapping time increases because there are more constraints to verify. It takes many iterations before finding the good t_b . Globally the increase of waiting time is due to class 3 applications (dedicated resources applications) which have to wait a long time before being alone on the machines. The economical or political criteria, which define important jobs, could depreciate the global utilization of resources. In case 1, there are 21 rejects of applications belonging to class 1 (deadline). In case 3, there are 154 rejects of applications belonging to class 1 (deadline). In those cases the machines are full. Because of the constraints of deadline, the tasks can not start later like in classical batch schedulers. Applications of class 2 can be refused but a mapping has always been found. Applications of class 3 and 4 can never be refused, they will only be slowed down.

5 Conclusion

This article presents a scheduling algorithm with quality of service usable in distributed systems like clusters or grids. It is implemented in AROMA : a resource management system used in ASP model. The validation shows that the proposed algorithm is better than NQS, nevertheless, the comparison is difficult because NQS does not implement classes of service and has not access to the same information.

The algorithm always respects the quality of service required for accepted applications. When different applications are mixed, the mapping and the waiting times are more important than when there are only best effort applications. The performances of the algorithm are still very good (mapping time inferior to 80 ms). The communication weight are introduced to favour the execution of an application in the same network area.

Future work will be to improve the notion of communication model and its use into the application model. The theoretical complexity of the mapping algorithm will be studied. The main difficulty will be to explore the complexity of the estimation of the end of a task. More comparisons with other mapping algorithms will be done. Real execution of a set of jobs during many weeks will be done on a small grid over 3 different sites.

References

1. Warren Smith, Ian Foster and Valerie Taylor. *Scheduling with advanced reservations*. Proceeding of the IPDPS Conference, May 2000.
2. Dror G. Feitelson and Morris A. Jette. *Improved utilization and responsiveness with gang scheduling*. Proceeding JSSPP 1997 : 238-261. Job scheduling strategies for parallel processing, IPPS'97 workshop, Geneva, Switserlang.
3. Dmitry Zotkin and Peter J. Keleher. *Job-length estimation and performance in backfilling schedulers*. 8th Intl Symp. High Performance Distributed Comput., august 1999.
4. P.Bacquet, O.Brun, J.M.Garcia, T.Monteil, P.Pascal, S.Richard. *Telecommunication network modeling and planning tool on ASP clusters*. Proceedings of the International Conference on Computational Science (ICCS'2003) Melbourne, Australia, June 2-4, 2003.
5. Atsuko Takefusa, Satoshi Matsuoka, Henri Casanova, Francine Berman. *A study of deadline scheduling for client-server systems on the computational grid*. Proceedings of the Tenth IEEE Symposium on High Performance Distributed Computing (HPDC10) San Francisco, California, August 7-9, 2001.
6. Rajkumar Buyya. *Economic-based distributed resource management and scheduling for grid computing*. Thesis, April 2002.
7. T.L. Casavant and J.G. Kuhl. *Effects of Response and Stability on scheduling in distributed computing systems*. IEEE Transactions on software engineering, vol. 14, No 11, pp. 1578-1588, november 1988.
8. Y.C. Chow, W.H. Kohler. *Models for dynamic load balancing in a heterogeneous multiple processor system*. IEEE Transactions on computers, vol. c-28, No 5, pp. 354-361, 1979
9. C.Y. Lee. *Parallel machines scheduling with non simultaneous machine available time*. Discrete Applied Mathematic North-Holland 30, pp 53-61, 1991.
10. F. Bonomi and A. Kumar. *Adaptative optimal load balancing in a non homogeneous multiserver system with a central job scheduler*. IEEE Transactions on computers, vol. 39, No 10, pp. 1232-1250, october 1990.
11. Henri Casanova, Arnaud Legrand and Loris Marchal *Scheduling Distributed Applications: the SimGrid Simulation Framework*. Proceedings of the third IEEE International Symposium on Cluster Computing and the Grid (CCGrid'03).
12. Patricia Pascal and Thierry Monteil. *Influence of Deterministic Customers in Time Sharing Scheduler*. ACM Operating Systems Review, 37(1):34-45, January 2003.
13. <http://www.cs.huji.ac.il/labs/parallel/workload/logs.htmlsdscsp2>
14. <http://joblog.npaci.edu/>
15. <http://www.cs.huji.ac.il/labs/parallel/workload/>
16. B. Kingsbury. *The network queuing system*. 16 May 1998. <http://pom.ucsf.edu/srp/batch/sterling/READMEFIRST.txt>.