

# MetaLoRaS: A Predictable MetaScheduler for Non-dedicated Multiclusters<sup>\*</sup>

J. Ll L rida<sup>1</sup>, F. Solsona<sup>1</sup>, F. Gin <sup>1</sup>, M. Hanzich<sup>2</sup>, P. Hern andez<sup>2</sup>, and E. Luque<sup>2</sup>

<sup>1</sup> Departamento de Inform tica e Ingenier a Industrial, Universitat de Lleida, Spain  
{jlerida, francesc, sisco}@diei.udl.es

<sup>2</sup> Departamento de Arquitectura y Sistemas Operativos, Universitat Aut noma de Barcelona, Spain  
mauricio@aomail.uab.es, {porfidio.hernandez, emilio.luque}@uab.es

**Abstract.** The main aim of this work is to take advantage of the computer resources of a single organization or Multicluster system to execute distributed applications efficiently without excessively damaging the local users. In doing so we propose a Metascheduler environment named MetaLoRaS with a 2-level hierarchical architecture for a non-dedicated Multicluster with job prediction capabilities.

Among other Metaschedulers, the non-dedicated feature and an efficient prediction system are the most distinctive characteristics of MetaLoRaS. Another important contribution is the effective cluster selection mechanism, based on the prediction system.

In this article, we show how the hierarchical architecture and simulation mechanism are the best choices if we want to obtain an accurate prediction system and, at the same time, the best turnaround times for distributed jobs without damaging local user performance.

## 1 Introduction

In this paper, we consider the issue of metascheduling jobs across a Multicluster. A Multicluster system has a network topology made up of interconnected clusters and limited to a campus- or organization-wide network with a common domain name, while a traditional grid is national or even global in scale [6]. Instead of grid computing, Multicluster systems allow the construction of efficient prediction systems based on the availability and reliability of their constituent resources.

Several studies [2] have revealed that a high percentage of computing resources in a Network Of Workstations (NOW/Cluster) are idle. The possibility of using this computing power to execute parallel jobs with a performance equivalent to a Massively Parallel Processor (MPP) and without perturbing the performance of the local user applications on each workstation has led to a proposal for new resource management environments ([7,11]).

With the aim of taking advantage of these idle computing resources available across the cluster, in previous work [4,5], we developed a cluster scheduling system named LoRaS (Long Range Scheduler). The most important feature of LoRaS is its efficient

---

<sup>\*</sup> This work was supported by the MEyC-Spain under contract TIN 2004-03388.

turnaround predictor for the parallel jobs, which allows the most efficient space sharing scheduling policy to be selected depending on the cluster state. In this article, we present MetaLoRaS, an extension of LoRaS, which provides the ability for a non-dedicated Multicluster to act as a Metascheduler. MetaLoRaS deals with each cluster by means of the LoRaS scheduling system requesting the cluster state and the prediction of the job turnaround time.

MetaLoRaS has a two-level hierarchical architecture. This choice is based on the work performed in [3], where an evaluation of the interaction between Metaschedulers and local schedulers (in each cluster) is presented. One important conclusion they presented was that among all the strategies considered, the best is to have one scheduler in each cluster, and to drop the requirement for scheduling single-component jobs to the cluster (or bottom) level.

The major scheduling decision of the top level of MetaLoRaS is to select the cluster to which the parallel jobs are delivered for execution. MetaLoRaS decisions are based on minimizing the execution time of the jobs without adding an excessive overhead to the local tasks. Parallel job execution times in each cluster are predicted in their respective LoRaS system, the bottom level of MetaLoRaS.

Most cluster selection policies used by Metaschedulers in the literature, employ traditional First Fit, Best Fit, or Greedy scheduling, or some variant on these ([13,10,9,3]). Moreover, an integral space sharing scheduling with load redistribution between clusters was used in [1]. A more sophisticated genetic algorithm was presented in [8]. An alternative technique to distribute the jobs between the clusters was proposed in [12] and [9]. They proposed distributing each job to the least loaded clusters in a redundant manner, thus improving the system utilization and reducing average job slowdown and turnaround time. However, the overhead involved in implementing this was found to be a problem.

Although these above-explained environments can give acceptable performance, they are not envisaged for non-dedicated environments and they cannot be used to predict job execution times in non-dedicated Multiclusters, the most important aims of this work.

In this article our metascheduling system was tested in a real Multicluster made up of three clusters. The results obtained demonstrate the good behavior of our Metascheduling mechanism and its applicability to Multicluster systems. Special attention was paid to the main contribution of this article, the cluster scheduling policy, that is, the algorithm which assigns jobs to clusters.

The remainder of this paper is outlined as follows. In Section 2, the Multicluster Scheduling system (MetaLoRaS) is presented. The efficiency measurements of MetaLoRaS are performed in Section 3. Finally, the main conclusions and future work are explained in Section 4.

## 2 MetaLoRaS

MetaLoRaS is based on a previously developed cluster scheduling system named LoRaS. LoRaS is a space sharing scheduler which allows a large number of scheduling facilities to be defined.

In this section the main features of LoRaS are introduced. Next, the Multiccluster architecture is extensively explained. In doing so, the main components of MetaLoRaS and their functionality are defined. Finally, the Metascheduling mechanism, consisting of selecting the cluster where the jobs are mapped, is explained separately.

## 2.1 LoRaS

LoRaS ([4,5]) is basically a Queue Manager developed in order to provide a space sharing scheduling facility for non-dedicated cluster systems. As can be seen in Fig. 1, the most important components of LoRaS are the Input Queue, the Predictor and the Scheduler.

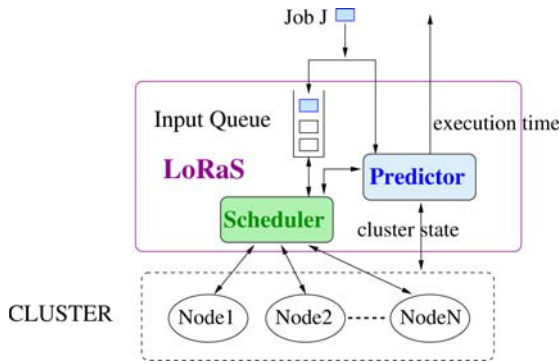


Fig. 1. LoRaS

The jobs in the Input Queue are ordered in a FCFS manner. The mapping policy of LoRaS selects the nodes for executing a parallel application considering the level of resource usage throughout the cluster. Thus, it does not overload any node without damaging the local user interactiveness.

LoRaS provides a means for predicting the execution time of parallel jobs in each cluster. This prediction is performed by simulation. For our purpose, the simulator is the most important component of LoRaS, because it is the basis for the Metascheduling mechanism proposed in this article.

In order to deal with this estimation, LoRaS uses information about the parallel applications and the cluster state. Each parallel application is executed in isolation and the following information is gathered: number of requested nodes, execution time and percentage of used CPU and Memory. The cluster state is modeled according to the following information about each cluster node: average load, memory occupancy, local user activity and the Multi-Programming Level (MPL), which is the number of parallel applications running in the node. If LoRaS sets the maximum MPL as  $N$ , no more than  $N$  parallel tasks can be executed in each cluster node.

LoRaS also deals with user activity. The nodes where there are local user activity, the MPL is adjusted accordingly. This way the local user jobs are also considered when

assigning new parallel tasks. The Metascheduler proposed in the next section improves the treatment of the local applications significantly.

As was demonstrated in [4] and [5], taking into account the information described above, the predictor can obtain accurate estimates of the turnaround time of the distributed jobs. In fact, the turnaround deviation achieved by the predictor system is always below 20%.

2.2 MetaLoRaS Architecture

First of all, we must decide whether the Multicluster system architecture will be centralized, hierarchical or distributed. Multiclusters are usually based on high-speed local networks, and therefore a centralized or hierarchical scheme should be suitable.

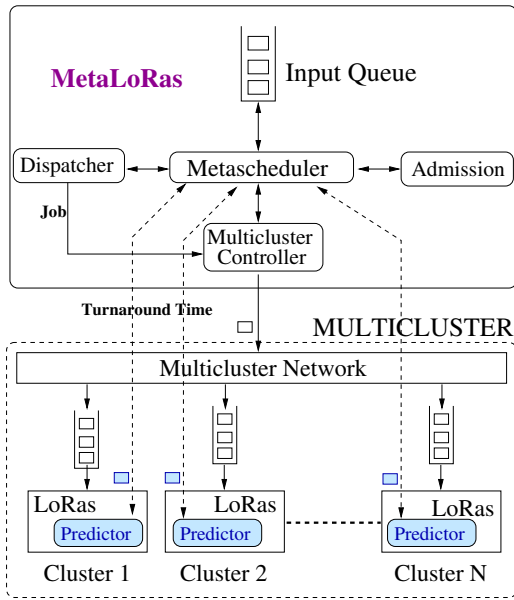


Fig. 2. MetaLoRaS

We are interested in extending LoRaS to make a Metascheduler system oriented to non-dedicated Multiclusters. The final result is MetaLoRaS. Unlike a centralized scheme, where the scheduling process should be performed in the Metascheduler node, we chose the hierarchical model. This decision was based on the attempt to split metascheduling decisions between the overall cluster or more precisely between the LoRaS front-end nodes, thus reducing the added scheduling overhead in the system. Furthermore, as was pointed out in [3], the best results were obtained when there was one scheduler in each cluster.

MetaLoRaS is made up of five components (see Fig. 2): the Input Queue (a queuing system), the Multicluster scheduler (named Metascheduler), the Admission system, the Dispatcher and the Multicluster Controller.

It is reasonably expected that users outside the organization should not have direct access to the local resources. Instead, they should use the Multicluster through one entry point, the Input Queue. Note that the Input Queue should comprise more than one queue. Doing so, the jobs could be classified, for example, on the basis of their scheduling priority. Nevertheless, for simplicity reasons, only one Input Queue is considered in this article.

*MetaScheduler*, the MetaLoRaS scheduler, is the Multicluster scheduling system. It is responsible for selecting the next job to be executed from the Input Queue, and also the cluster where this job will be executed. *Metascheduler* is explained in depth in section 2.3.

The *Admission System* is responsible for admitting new jobs into the system. This module will accept the new job whenever its required resources fit on at least one cluster. If not, the job is discarded. The specified resources are the number of processors and the size of Memory.

The MetaLoRaS maintains an updated list of the state of each cluster, so it can determine in real time whether the first job in the Input Queue can be executed. This information is obtained from the Multicluster Controller system. The *Multicluster Controller* collects real time information about the state of each cluster. If an event occurs in one cluster (job start, finish), the Multicluster Controller is notified of such a change. It is responsible for updating the database with the information from the different clusters. In doing so, it collects the changes performed in each cluster. The LoRaS system is responsible for notifying the Multicluster Controller about the cluster state changes.

The *Dispatcher* has the ability to send the job to the Input Queue of the LoRaS frontend, selected by MetaLoRaS. The jobs to be executed can be launched in both PVM and MPI formats.

One important issue to be taken into account is the Multicluster Network topology. The nodes are usually grouped into clusters because they share an internal network. A Multicluster has a distinctive architectural feature: the internal networks of the clusters are bridged together through dedicated links. This configuration increases the complexity of effectively managing both computing and networking resources.

Finally, we suppose that the clusters making up the Multicluster are accessible from the node where the MetaLoRaS is located. Also, we suppose that the clusters are shared between local user applications and the parallel jobs.

## 2.3 MetaLoRaS Scheduling System

In this section the functioning of MetaScheduler, the Multicluster scheduling system (see Algorithm 1) is explained. MetaScheduler is responsible for selecting jobs from the Input Queue and mapping them into clusters, controlled by the LoRaS system.

The next job ( $J$ ) from the Input Queue to be executed is selected in a FIFO manner. In this article, as we are interested in presenting an efficient prediction system, only the FIFO policy is considered because optimal prediction accuracy is obtained when using this policy [5]. If the selected job  $J$  cannot be executed in the Multicluster, i.e. the resources that the job requires do not fit in one cluster, the admission system will reject the job.

Each job must provide the number of processors and their respective CPU and Memory requirements. As mentioned in the previous section, they are obtained by the LoRaS when executing in isolation.

---

**Algorithm 1.** Cluster Selection Policy.

---

**Step 1** Choose in a FIFO manner the next job ( $J$ ) from the Input Queue

**Step 2** If there are not enough resources in every Cluster, reject job  $J$  and **Go to** Step 1

**Step 3** Obtain the Cluster  $C$  to map job  $J$  according to the Prediction Algorithm (Algorithm 2).

**Step 4** Dispatch Job  $J$  to Cluster  $C$

**Step 5** **Go to** Step 1

---

If the selected job  $J$  is not rejected (fits in almost one cluster), it will be mapped into one cluster by means of the Cluster Selection Policy (see algorithm 1). Jobs are selected in a FIFO manner and the considered resources are the Memory occupancy and the maximum fixed MPL of the cluster nodes. This information is maintained by the Multicluster Controller system. For this reason, each cluster must provide information to the Multicluster Controller when jobs complete and free processors.

The proposed policy is based on simulation. The job execution is simulated in each LoRaS simulator, residing in the front-end node of each cluster. With this model, the clusters do not perform any scheduling decision, but are only responsible for dispatching the jobs that are supplied by the MetaLoRaS. The *Prediction Algorithm* (Algorithm 2) summarizes the steps performed by the simulation process in obtaining the cluster to map the job  $J$ .

Next, the *Prediction Algorithm* is explained. The turnaround time of a job  $J$  is obtained by simulation in all the clusters. In doing so, the cluster state, the waiting jobs in the LoRaS Input Queue and the local user activity is taken into account. To select the cluster, the MetaLoRaS is based on the Pondered Turnaround Time metric ( $PTT(J)$ ).

Let  $C$  be a Multicluster made up of “ $n$ ” clusters. The Pondered Turnaround Time for a cluster  $k \leq n$  ( $PTT_k$ ) and a job  $J$ , is defined as:

$$PTT_k(J) = TR_k(J) \left( W \frac{TR_k(J)}{MaxTR(J)} + (1 - W) \frac{LT_k}{MaxLT} \right), \quad (1)$$

where  $TR_k(J)$  is the predicted turnaround time obtained in cluster  $k$ ,  $MaxTR(J)$  is the maximum predicted turnaround time across the Multicluster,  $LT_k$  is the number of nodes assigned to job  $J$  in cluster  $k$  with local user activity and finally,  $MaxLT$  is the maximum number of nodes assigned to job  $J$  with local activity across the Multicluster. To prevent a division by 0,  $MaxLT$  is set to 1 when there is no local activity at all ( $MaxLT=0$ ).  $W$  (a real value in the range  $[0..1]$ ) is the weight assigned through the formula to the turnaround time.  $1 - W$  will then be the weight assigned to the resource occupancy of user applications.

We define the Pondered Turnaround Time of a job  $J$  as the minimum  $PTT_k \forall k = 1..n$ . This is:

$$PTT(J) = \min_{k=1}^n \{PTT_k(J)\} \quad (2)$$

Our Metascheduling system will assign the job  $J$  the cluster which gives the minimum turnaround time and that disturbs the local user applications as little as possible. According to equations 1 and 2, Cluster  $k$  with  $PTT(J) = \min_{k=1}^n \{PTT_k(J)\}$  will be selected.

---

**Algorithm 2.** Prediction Algorithm.

---

**Step 1** Metascheduler obtains the list of clusters ( $L_c$ ) where the job can be executed.

**Step 2** For each cluster  $C \in L_c$ , LoRaS obtains its cluster state.

**Step 3** Metascheduler request the turnaround time of Job  $J$  in each cluster  $k$  where the job can be executed ( $TR_k(J)$ ). Also request the local resource occupancy of each cluster ( $LT_k, \forall k = 1..n$ ).

**Step 4** Metascheduler computes  $PTT_k(J), \forall k = 1..n$ .

**Step 5** Metascheduler obtains the cluster  $C$  where  $PTT(J) = \min_{k=1}^n \{PTT_k(J)\}$  is reached.

---

By varying  $W$  we can obtain different scheduling capabilities. If the local user applications must not be delayed by the parallel jobs, we will set  $W \simeq 0$ . Instead, If the major performance metric to be taken into account is the turnaround time of parallel applications, a  $W \simeq 1$  should be used. For a value of  $W \simeq 0.5$ , equal importance will have both local user intrusion and parallel turnaround times.

Note that there can also be different  $W$  values for different cluster systems. This should be used for example to distinguish between clusters where the user applications cannot be damaged at all ( $W \simeq 0$ ).

### 3 Experimentation

In order to carry out the experimentation process, we need the user and parallel workloads. The user workload is represented by a synthetic benchmark (named *local\_bench*) which can emulate the usage of 3 different resources: CPU, Memory and Network traffic. The use of these resources was parametrized in a real way. According to the values obtained by collecting for a couple of weeks the user activity in an open laboratory, *local\_bench* has been modeled to use 15% CPU, 35% Memory and 0,5KB/sec LAN.

The parallel workload was a list of 60 NAS parallel jobs (CG, IS, MG, BT), with a size of 2, 4 or 5 tasks, that reached the system following a Poisson distribution. These jobs were merged so that the entire parallel workload had a balanced requirement for computation and communication. It is important to mention that the parallel MultiProgramming Level (MPL) reached for the workload depends on the system state at each moment, but cannot exceed an  $MPL = 4$ .

In this experimentation we are interested in knowing the benefits of the new *Cluster Selection Policy* (named PTT) based on the state of the Multicuster. In order to be able to show the benefits of this policy, a comparison was done between different policies with different values of local activity. The policies compared are: FFIT (First Fit), RR (Round Robin), TAT (Turnaround Time) and PTT (Pondered Turnaround Time). The first two are classic policies, whereas TAT and PTT use the LoRaS prediction system. TAT assigns jobs in the cluster with the minimum turnaround time. In PTT (the policy

presented in this article), not only is the minimum turnaround time taken into account (as in TAT), but also low user intrusion is considered, as was expressed in Formula 2.

The whole system was evaluated in a Multicluster made up of three non-dedicated clusters with the same number of nodes. The 3 clusters, named Cluster1, Cluster2 and Cluster3 each had 5. Each cluster comprised 3GHz workstations with 1GB of RAM and 2048 KB of cache, interconnected by a 1Gigabit network.

### 3.1 Metrics

In order to be able to compare the different policies three different metrics, *Turnaround*, *Makespan* and *Slowdown*, were used.

The *Turnaround* measures the mean elapsed time in executing one job of the parallel workload. This metric allows us to measure the efficiency of the different Cluster Selection Policies evaluated from the parallel application point of view.

*Makespan*, on the other hand, allows us to analyze the global performance of the overall system. It gives the elapsed time from the start of the first job to the end of the last job. This metric relates the performance of the system in executing the overall parallel workload. It is considered a system centric metric.

To quantify the intrusion introduced by the execution of the parallel workload in the user applications, we used a benchmark, called *local\_bench*, that measures the averaged latency of a set of system calls. The *Slowdown* (*SL*) of local user applications is defined as follows:

$$SL = \frac{AvLat_{non-dedicated}}{AvLat_{dedicated}} \quad (3)$$

where  $AvLat_{non-dedicated}$  is the averaged latency value in executing *local\_bench* with the parallel workload (non-dedicated environment) and  $AvLat_{dedicated}$  is the averaged latency of the same benchmark when executed in isolation (dedicated environment).

### 3.2 Performance Comparison

In the Figure 3 we show the averaged Turnaround time of the parallel workload by varying the local activity. 0% means that there are no local users in any cluster; 30% one cluster is entirely occupied by user applications; 60% two clusters are occupied and finally 100%, the overall Multicluster has user activity.

As can be seen in the figure, the FFIT policy is the worst. In this case, almost all the parallel jobs were assigned to the same cluster (the first one) because there were enough resources. Only when there were not enough processors to execute the job was another cluster selected.

The RR policy behaved better than FFIT in all the situations, because it distributed better the jobs. As Figure 3 shows, the policies TAT and PTT mapped even better because the assignment is based on the prediction of the turnaround time. This improvement also demonstrates the effectiveness of the prediction module of LoRaS.

As we can see, the behavior of the PTT model depends on the value of  $W$  and the user workload. PTT with  $W=1$  gave similar performance to TAT in all the cases. As Formula 2 shows, this result was expected because in both cases only the turnaround time is considered.



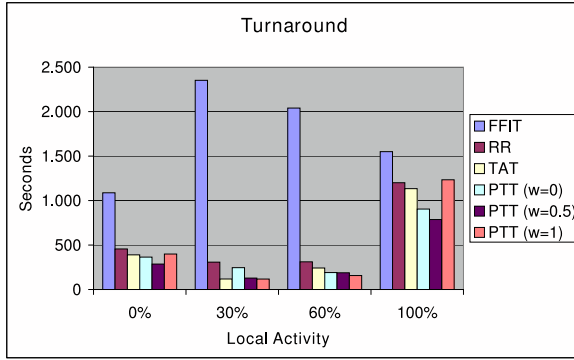


Fig. 3. Turnaround

For the  $W=0$  case, the PTT model discards the clusters with local activity. When there are various clusters without user activity, the policy selects the first cluster with enough resources. In many situations the assignments are performed in the same cluster, penalizing the Turnaround metric, as passed in the 30% case.

The best performance results are obtained for PTT with  $W=0.5$  and  $W=1$ . This is due to the fact that in both cases the turnaround time is considered, avoiding the commented problems of  $W=0$ . The average of the performance results is similar in both cases.

PTT increased the performance of RR and FFIT in all the situations. Also, on average, the PTT model with  $W=0.5$  and  $W=1$  improved TAT in 5%.

The Figure 4 shows the Makespan results. This metric measures the overall system performance with respect to the applied policy. As was expected, the obtained results are very similar to the turnaround. This fact corroborates the good behavior of our experimental environment and the implemented policies.

With regard to the Makespan metric, PTT is also the best policy. Obtaining, in average, a gain with respect to TAT of 2% for PTT  $W=0$  and 8% for PTT with  $W=0.5$  and  $W=1$ .

### 3.3 User Slowdown

In this section, the influence of the different policies on the user applications is evaluated.

Figure 5 shows the Slowdown obtained when executing *local\_bench* in the same situations as in the previous section. As can be seen, the intrusion does not increase with the local activity. This demonstrates that LoRaS deals with the local user well, by balancing the maximum MPL allowed in the cluster with the local activity. The greater the local activity, the smaller the MPL and this causes that the latency to stay stable.

In general, a great negative impact on the local task performance was produced by the FFIT model. In this case, only taking into account the fitting of the jobs in the assignment increased the intrusion in the user applications. The RR model also does so. With a simple distribution of the load it also caused a considerable overhead in the user workload.

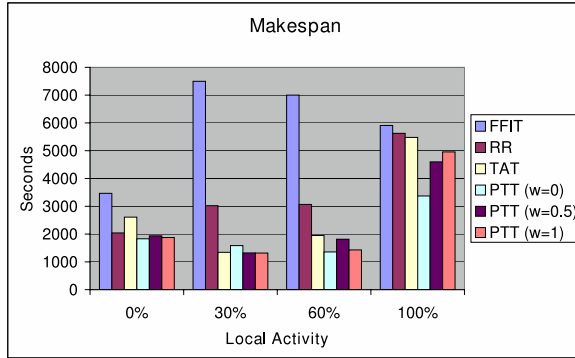


Fig. 4. Makespan

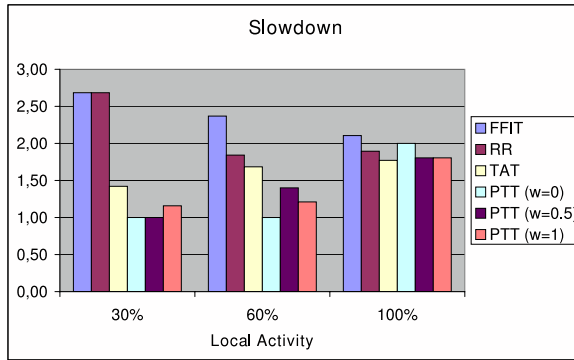


Fig. 5. Local Intrusion

We can observe how PTT maintains the smallest latency in almost all the cases. The exception occurred in the case of  $W=0$  and 100% user activity. This was produced because this model tries to decide only according to the user activity. As all the clusters have the same user workload, this model decides in arbitrarily.

We have verified that in average, the PTT model with  $W=0$  obtained the smallest latency. However, the TAT model improved performance of parallel applications with respect to PTT with  $W=0$  in 11%. We have also verified that the best performance is obtained for PTT with  $W=0.5$  and  $W=1$ , with a gain with respect to TAT of 5%. Furthermore, PTT with  $W=1$  obtained smaller latency values than PTT with  $W=0.5$ .

## 4 Conclusions and Future Work

In this article we have presented MetaLoRaS, a Metascheduler system based on simulation with a two-level hierarchical architecture for a Multicluster environment. Compared with traditional policies in the literature, the simulation model we propose gave the best cluster selection scheduling because it can predict the turnaround time of an

application in a particular cluster system. This way, no redundant job assignments must be implemented to obtain the best result (as in [12]).

The proposed Metascheduler gives good parallel performance without excessively damaging the user workload. This is one of the most important aim of this work. Our Metascheduler is very efficient when it is based only on the prediction mechanism. It reduced the Makespan of the parallel workload by 23% and the turnaround by 25% with respect to the classic policies (First Fit and Round Robin).

The proposed new policy, PTT, based on the prediction mechanism and the user activity allows the intrusion of user applications to be reduced by an average of 9%. Varying the  $W$  factor, we obtain different performance results depending on the user activity. For the  $W=0$  case, we obtained better *Slowdown* results. Instead, for the  $W=0.5$  model, we obtained better average turnaround and Makespan values.  $W=1$  obtained gains similar to  $W=0.5$  but reducing the user intrusion.  $W=1$  is the best mode if we want to balance parallel against user application performance. To improve the overall system performance with low user intrusion,  $W=0$  should be chosen. Instead,  $W=0.5$  is the best model if the goal is to improve the parallel application performance.

Future work is directed towards improving the scheduling model when the user activity is similar across the Multicluster. As our model only considers the presence of user activity, it do not select efficiently the clusters to map the jobs when the clusters are evenly loaded. We need additional information, such as CPU or Memory occupancy to distinguish the least loaded cluster to perform more efficient decisions. We want also consider the workload changes in advance, making load-redistribution (as in [1,9]) unnecessary.

The number of jobs waiting in the Input Queue of the LoRaS systems reduce prediction accuracy. To correct this problem new proposals in the PTT policy should be made.

## References

1. J. Abawajy and S. Dandamudi. Parallel Job Scheduling on Multicluster Computing Systems. In Proceedings of the IEEE International Conference on Cluster Computing (CLUSTER'03). 2003.
2. A. Acharya and S. Setia. Availability and utility of idle memory in workstation clusters. In Proceedings of the ACM SIGM./PERF'99, pp. 35–46. 1999.
3. A. Bucur and D. Epema. Local versus Global Schedulers with Processor Co-allocation in Multicluster Systems. JSSPP 2002, LNCS 2537, pp. 184–204. 2002.
4. M. Hanzich, F. Giné, P. Hernández, F. Solsona and E. Luque. Coscheduling and Multiprogramming level in a non-dedicated cluster. LNCS vol. 3241, pp. 327–336. 2004.
5. M. Hanzich, F. Giné, P. Hernández, F. Solsona and E. Luque. Cisne: A new integral approach for scheduling parallel applications on non-dedicated clusters. LNCS vol. 3648, pp. 220–230. 2005.
6. W. Jones, W. Ligon, L. Pang. Characterization of Bandwidth-Aware Meta-Schedulers for Co-Allocating Jobs Across Multiple Clusters. The Journal of Supercomputing, vol. 34, pp. 135–163. 2005.
7. M. Litzkow, M. Livny, and M. Mutka. Condor- a hunter of idle workstations. 8th Int'l Conference of Distributed Computing Systems. 1988.

8. A. Bose, B. Wickman and C. Wood. MARS: a metascheduler for distributed resources in campus grids. Proceedings. Fifth IEEE/ACM International Workshop on Grid Computing. 2004.
9. G. Sabin, R. Kettimuthu, A. Rajan and P. Sadayappan. Scheduling of Parallel Jobs in a Heterogeneous Multi-site Environment. JSSPP 2003, LNCS 2862, pp. 87–104. 2003.
10. J. Santos, G.D. van Albada, B.A. Nazief, P.M. Sloot. Hierarchical Job Scheduling for Clusters of Workstations. Proceedings of the sixth annual conference of the Advanced School for Computing and Imaging (ASCI 2000), pp. 99-105. 2000.
11. E. Shmueli and D. G. Feitelson. Backlling with lookahead to optimize the performance of parallel job scheduling. Job Scheduling Strategies for Parallel Processing, LNCS, 2862:228251. 2003.
12. V. Subramani, R. Kettimuthu, S. Srinivasan and P. Sadayappan. Distributed Job Scheduling on Computational Grids using Multiple Simultaneous Requests. Proceedings of the 11<sup>th</sup> IEEE International Symposium on High Performance Distributed Computing. 2002.
13. M. Xu, "Effective Metacomputing using LSF MultiCluster," ccgrid, p. 100, 1st International Symposium on Cluster Computing and the Grid, 2001.