

Topology-aware resource management for HPC applications

Yiannis Georgiou
ATOS/Bull
Grenoble, France
yiannis.georgiou@atos.net

Emmanuel Jeannot
Inria Bordeaux Sud-Ouest
Talence, France
emmanuel.jeannot@inria.fr

Adèle Villiermet
Inria Bordeaux Sud-Ouest
Talence, France
adele.villiermet@inria.fr

Guillaume Mercier
Bordeaux INP
Talence, France
guillaume.mercier@bordeaux-inp.fr

ABSTRACT

The Resource and Job Management System (RJMS) is a crucial system software part of the HPC stack. It is responsible for efficiently delivering computing power to applications in supercomputing environments. Its main intelligence relies on resource selection techniques to find the most adapted resources to schedule the users' jobs. Improper resource selection operations may lead to poor performance executions and global system utilization along with an increase of the system fragmentation and jobs starvation. These phenomena play a role in the increase of the platforms' total cost of ownership and should be minimized. This paper introduces a new method that takes into account the topology of the machine and the application characteristics to determine the best choice among the available nodes of the platform based upon their position within the network and taking into account the applications communication pattern. To validate our approach, we integrate this algorithm as a plugin for SLURM, a popular and widespread HPC resource and job management system (RJMS). We assess our plugin with different optimization schemes by comparing with the default topology-aware SLURM algorithm using both emulation and simulation of a large-scale platform, and by carrying out experiments in a real cluster. We show that transparently taking into account the job communication pattern and the topology allows for relevant performance gains.

CCS Concepts

•Information systems → Computing platforms;

Keywords

resource management; job allocation; topology-aware placement; scheduling; SLURM

1. INTRODUCTION

Computer science is more than ever a cornerstone of scientific development, as more and more scientific fields resort to simulations in order to help refine the theories or conduct experiments that cannot be carried out in reality because their scale or their cost are prohibitive. Currently, such computing power can be delivered only by parallel architectures. Larger and larger machines are being built around the world, and being able to display such a machine has become a challenge for states and nations, scientifically as well as politically.

However, harnessing the power of a large parallel computer is no easy task, due to several factors. First, this type of computer features usually a huge amount of computing nodes and this scale has to be taken into account when developing applications. Then, the nodes architecture has become more and more complex, as the number of cores per node is in constant increase from one generation of CPU to the next. The memory hierarchy becomes also more complex, as various levels of cache are now available and the rise of MCDRAM or NVRAM will make things even more complicated in the future. Indeed, an efficient exploitation of all these types of memories is possible only if the application developer takes it into account.

One way of dealing with such complexity would be to consider the application behavior (e.g. its communication pattern, or its memory accesses pattern) and to deploy it on the computer accordingly. To this end, the most widespread technique is to determine the list of cores on which the application has to be run on, and then to bind the processes on these cores so as to minimize/maximize a predetermined criterion (a.k.a. a metric). Such a technique has already been used and investigated to improve the performance of parallel applications [14].

However, a large parallel machine is often shared by many users running their applications concurrently. In such a case, an application execution will depend on its nodes allocation, as determined by the Resource and Job Management System (RJMS). Most of the time RJMS work in a best-effort fashion, which can lead to suboptimal allocations. That is, such allocations might be able to fulfill an application requirements in sheer terms of resources (number of CPUs, amount of memory) but might also fail to provide an environment tailored for an optimized execution. For instance, if

the application processes communicate a lot between themselves, a set of nodes physically allocated apart from the rest might degrade performance severely. Furthermore, even if the given allocation is contiguous, taking into account process affinity leads to even better performance.

As a consequence, our goal is to apply to resource management the same technique that has proved its efficiency for applications deployment and execution, that is, taking into account an application's behaviour in the process of reserving and allocating the needed resources (computing nodes). This means more criteria to be used and considered by the RJMS when a user submits its request to the system. Actually, taking in account an application behaviour when allocating nodes pushes even further the idea of using an application information to improve its execution.

In this paper, we shall detail the improvements we made to an existing RJMS in order to enable it to select the most suitable set of nodes for a given parallel application. To this end, we extend our TREEMATCH algorithm and integrate it in the SLURM software to improve its ability to match the resources to the actual application communication pattern. This paper is organized as follows: Section 2 gives an overview of the context and background of this work. Section 3 introduces all the software leveraged by this work before giving more technical insights about our topology-aware job allocation policy. Then Section 4 shows and discusses the results obtained while related work are listed in Section 5. Finally, Section 6 concludes this paper.

2. ISSUES OF RESOURCE ALLOCATION IN PARALLEL COMPUTERS

2.1 The Sharing of Resources

A large parallel computer is to some extent a tool that has to be exploited and used. A reason as why these computers increase in size and scale stems from the fact that some applications grow accordingly. Therefore, an adequate platform has to match these needs. However, a substantial part of the time, this large platform not only works in a time-sharing mode, but also in a space-sharing mode. Indeed, in order to exploit the hardware in a satisfactory way, several users share it, leading to a potentially very large number of users. An interactive access is therefore out of the question. To this end, the users have to submit their requests in terms of resources to a system called the Resource and Job Manager System (sometimes called a Batch Scheduler for short). This system's goals are threefold: 1. to centralize and analyze all the received requests, 2. to allocate the most relevant type of resources (CPU, memory or network switches for instance) able to fulfill these demands and 3. to execute the application (a.k.a. the job) submitted by a user on the set of selected resources.

There are many and sometimes conflicting criteria that should be optimized by the RJMS. Then the question that pertains to this selection and allocation of resources is to choose one. For instance, one metric could be the system throughput, that is, the amount of jobs executed during a defined time step, whilst another could be the use (CPU load) of the system. All these metrics are relevant and which to use/optimize depends on a given point of view. That is, an administrator's point of view might diverge from a user's point of view. Indeed the users hardly possess a global

view of the system (in most of cases) as opposed to the administrators, hence the discrepancy.

2.2 Finding an Optimization Criterion

In this work, we focus on a metric relevant for users: the flow time (or turnaround time) that is, the time his/her job remains in the system. We believe it to be the most appealing one for a user seeking to gather results and get the outcome of his/her application as soon as possible. The question that now arises is how to speed up an application execution? Let us suppose that the developer has already optimized his/her application as much as it is possible. What are the means left to even speed things further up? One answer lies in the ecosystem of the application, that is, in the way the application is deployed and executed. In a previous work, unrelated to resource management and job scheduling, we showed that by taking into account an application behaviour when deploying it on the various processing entities (CPUs, cores, threads, etc.), it is possible to improve its global execution time [18, 15]. Actually, our goal is to improve the way an application accesses its data. This data locality can be improved in several ways, but we chose so far to use the communication pattern of the application, that is, an expression of the amount of bytes/messages exchanged by the application processes. Then, we try to match this pattern to the underlying architecture by following the principle that the more processes are communicating with the others, the closer the cores they should be bound to. This can be done by several techniques but usually involves process binding and rank reordering [19].

However, the execution still depends on the set of resources allocated to the application by the RJMS. Since no guarantee is given that this allocation will be compliant with the application communication pattern, some negative side effects may occur. For instance, a subset of nodes might be physically far from another subset, thus impacting the communication between processes belonging to each subset. As a consequence, an allocation that takes into account an application communication scheme leads to performance improvements. To that end, we consider a well-known and widespread RJMS called SLURM and design a new plugin based on the TREEMATCH algorithm. So far, TREEMATCH was used to compute a matching between the application processes and the physical cores available. Now, we use it to determine a nodes allocation before deploying the application.

Hence, to improve the flow time we aim at reducing the job execution time of the submitted application by improving its mapping.

2.3 A Motivating Example

The goal of this work is to apply mapping techniques (e.g. TREEMATCH) before the execution of the application processes and compare different approaches. We assume that the communication pattern of the application is known at submission time. Such a communication pattern can be gathered with application monitoring (see Section 3.1.2) or by analyzing the structure of the parallel algorithm (for instance if we are dealing with a stencil code we know which processes are communicating together and the amount of exchanged data). In any case, we assume that this communication pattern remains unchanged from one run to the other. It is not the case for all parallel applications but a large

Proc.	0-1	2-3	4-5	6-7
0-1	0	20	0	2000
2-3	20	0	1000	0
4-5	0	1000	0	10
6-7	2000	0	10	0

Table 1: Affinity matrix for 8 processes (4 groups of 2 processes each). Shows the amount of bytes/messages exchanged by the application processes

amount of applications comply to these models (for instance, dense linear applications and kernels, stencil codes, regular mesh partitioning based applications, etc.). When the communication pattern changes from one run to another, the proposed solution is not applicable and the user has to fall back to a standard allocation scheme: mixing the proposed topology-aware mapping with other types of mapping is totally acceptable.

Several possibilities are available. The most obvious one is to not use TREEMATCH at all and let the SLURM environment deal with the topology by itself. The second possibility is to apply TREEMATCH just before the job execution, once SLURM has selected the resources. Another possibility is to use TREEMATCH inside the selection mechanism of SLURM.

An example of the difference between these approaches is depicted by Fig. 1. Let us suppose that we have 6 nodes composed of two computing entities each. We assume that node n3 is not available as computing entities 6 and 7 are already used by another application, hence unavailable for a job allocation. Let us assume that a newly submitted job requests 4 nodes. For the sake of simplicity, we group processes in pairs (0-1, 2-3, etc.) and hence each pair of processes shall be assigned to one node. The affinity matrix is given in table 1.

If SLURM has to allocate resources for these 8 processes, it will look for the smallest number of switches able to fulfill the request. In this case, it will require to use the whole tree. Then, it will allocate processes from left to right inside nodes in a round-robin fashion. It will allocate nodes 0, 1, 2 and 4 for the job and then map processes onto the computing entities. We can see that such an allocation is rather costly communication-wise as groups of processes are spread onto the entities and no optimization is enforced in this regard. It is therefore possible to call TREEMATCH (see Section 3.1.3) to optimize the process mapping on these entities accordingly to the affinity matrix. By doing so, the resulting mapping is: group 0-1 on n0, group 6-7 on n1, group 2-3 on n2 and 4-5 on n4. This is the best possible solution once the resources have been allocated. However, group 2-3 communicates a lot with group 4-5. With such an allocation, all the communications will transit through the root of the topology, a costly solution in terms of hops. However, a better outcome is achievable if TREEMATCH performs the resource allocation. Given such a topology and the above affinity matrix, TREEMATCH will allocate group 0-1 on n0, group 6-7 on n1, group 2-3 on n4 and 4-5 on n5 since there are constraints on node n3.

In this case, all the communication between group 2-3 and group 4-5 will take only 2 hops instead of 4 and therefore the communication cost is even more reduced.

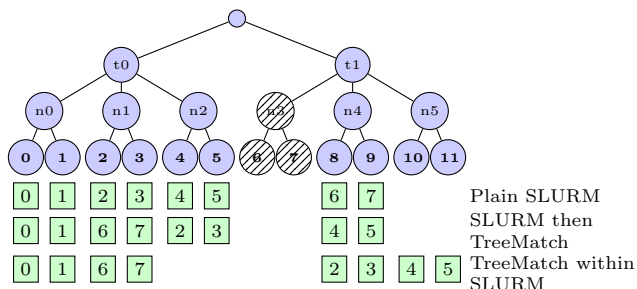


Figure 1: Tree topology of 6 nodes of 2 processing units with one unavailable nodes: n3

3. A TOPOLOGY-AWARE RESOURCE AND JOB MANAGEMENT SYSTEM

3.1 Software

In this section, we introduce with more details the various software elements that we use to implement the work described in this paper. First, we shall describe SLURM, our target RJMS. Then, we shall explain the method employed to gather information about the application communication scheme (a.k.a. our *affinity matrix*). Then, we give more specific information about the TREEMATCH algorithm and the constraints mapping extension we have implemented.

3.1.1 SLURM

We implement a new topology-aware placement algorithm within the open-source resource and job management system SLURM [30]. SLURM performs workload management on six of the ten most powerful computers in the world of the Top500 list¹ including the system ranked number one, Tianhe-2, which features 3,120,000 computing cores.

SLURM is specifically designed for the scalability requirements of state-of-the-art supercomputers. It is based upon a centralized server daemon, `slurmctld` also known as the controller, which communicates with client daemons `slurmd` running on each computing node. Users can request the controller for resources to execute interactive or batch applications, referred to as jobs. The controller dispatches the jobs on the available resources, whether full nodes or partial nodes, according to a configurable set of rules. The SLURM controller also features a modular architecture composed of plugins responsible for different actions and tasks such as: job prioritization, resources selection, task placement or accounting.

The resource selection process within SLURM takes place as part of the global job scheduling procedure. In particular, this procedure makes use of the `plugin/select`, which is responsible for allocating the computing resources to the jobs. Other plugins are used to facilitate and extend this procedure such as `plugin/topology` which takes into account the network topology of the cluster, the `plugin/gres` which can extend the allocation to different generic resources and the `plugin/task` which provides the isolation and possible binding of tasks on the resources.

There are various resource selection plugins within SLURM that can take into account the specificities of the underlying platforms' architecture such as `linear` and `cons_res`.

¹<http://top500.org/lists/2015/11/>

The `select/linear` plugin allows the allocation of complete nodes for jobs, using simple and scalable best-fit algorithms, however, the lower granularity of allocatable unit is the node which is quite limiting for new multicore and manycore architectures. The `select/cons_res` plugin is ideal for this type of architectures where nodes are viewed as collections of consumable resources (such as cores and memory). In this plugin, nodes can be used exclusively or in a shared mode where a job may allocate its own resources different than other jobs using the same node. The algorithms within the `cons_res` plugin are also scalable, featuring best-fit placement of jobs but they are more complex than `select/linear` since a finer granularity of allocatable resources is taken into account. One of the first version of the `select/cons_res` plugin is described in [2].

Our studies and developments as described in the following sections are based upon the `select/cons_res` plugin therefore we try to analyze a bit more some important internals of this plugin. The internal representation of resources and availabilities within SLURM is made using bitmap data structures. In the case of the linear plugin only a node bitmap is needed whereas in the case of the `cons_res` plugin, besides the node bitmap, a core bitmap is used to represent internal node resources availabilities. Within the `cons_res` plugin, the usage of node and core bitmaps is leveraged efficiently (e.g. kept separated in different contexts) in order to keep a high scalability for the selection algorithms. Another functionality of the `cons_res` plugin is the distribution of tasks within the allocated resources, which is an important feature for the optimal performance of parallel applications.

SLURM provides configuration options to make the resources selection network topology-aware through the activation of the topology plugin (`topology/tree` plugin). A particular file describing the network topology is needed and the job placement algorithms favor the choice of groups of nodes that are connected under the same network switch. The goal of the SLURM topology-aware placement algorithms is to minimize the number of switches used for the job and provide a best-fit selection of resources based on the network design. This feature becomes mandatory in the case of pruned butterfly networks where no direct communication exists between all the nodes. We use this plugin in our experiments. The scalability and efficiency of topology-aware resource selection of SLURM has been evaluated in [10].

Finally since the `cons_res` plugin deals with multi-core architectures the isolation and binding of tasks upon the used resources is an important feature to guarantee a minimal interference between jobs sharing nodes. This feature takes place through the usage of the `task/affinity` or the `task/cgroup` plugin which use linux kernel mechanisms such as cgroups and cpusets or APIs such as hwloc [5] in order to provide the described isolation and binding.

3.1.2 Application Monitoring

For this work we need to model an application communication scheme. The way communications occur describes the affinity between processes. For the affinity matrix, we gather the communication pattern thanks to a dynamic monitoring component we integrate in Open MPI as an MCA (Modular Component Architecture) framework called pml (point-to-point management layer). This component, when activated at launch time, monitors all the communications at the lowest level in the Open MPI stack (i.e. once collec-

tive communications have been decomposed into point-to-point operations). Therefore, as opposed to the standard MPI profiling interface (PMPI) approach where the MPI calls are intercepted, we monitor in our case the actual point-to-point communications that are issued by Open MPI, which is much more precise: for instance, we can see the tree used for aggregating values in a `MPI_Gather` call.

Internally, this component uses the low-level process ids and creates an associative array to convert sender and receiver ids into ranks in `MPI_COMM_WORLD`. At the end of the execution, each process dumps its local view into a file and a script aggregates all the local views at a given process to get the full communication matrix.

3.1.3 TreeMatch

TREEMATCH [15] [14], is a library for performing process placement based on the topology of the machine and the communication pattern of the application, for multicore, shared memory machines as well as distributed memory machines. It computes a permutation of the processes to the processors/cores in order to minimize the communication cost of the application.

To be more specific, it takes as input a tree topology (where the leaves stand for computing resources and internal nodes correspond to switches or cache levels) and a matrix describing the graph affinity between processes. The topology information is supplied either by the RJMS or by tools such as hwloc or netloc². A hierarchy is extracted from this graph so that it matches the hierarchy of the topology tree. The outcome is a mapping of the processes onto the computing resources. The objective function optimized by TREEMATCH is the Hop-Byte [31], that is, the number of hops weighted by the communication cost:

$$\text{Hop-Byte}(\sigma) = \sum_{1 \leq i < j \leq n} \omega(i, j) \times d(\sigma(i), \sigma(j))$$

where n is the number of processes to map, σ is the process permutation output produced by TREEMATCH (process i is mapped on computing resource $\sigma(i)$), $A = (\omega_{i,j})$ $1 \leq i \leq n$, $1 \leq j \leq n$ is the affinity matrix between these entities and hence $\omega(i, j)$ is the amount of data exchanged between process i and process j and $d(p_1, p_2)$ is the distance, in number of hops, between computing resources p_1 and p_2 . In a previous work [15], we have shown that minimizing this metric allows for application runtime reduction for tree-based topologies.

An important feature of TREEMATCH is that it only uses the structure of the tree and does not require a precise valuation of the speed of the links in the topology. Therefore, TREEMATCH does not require a performance assessment of the system on which the application is going to be executed. We believe this to be a strong advantage, as gathering such information is error-prone, might be incomplete and subject to inaccuracy.

In order to tackle the fact that not all resources are available for mapping we enhance TREEMATCH from [15] to take constraints into account. When not all leaves are available for mapping (because some of them are already used by other applications), it is possible to restrict the leaves onto which processes can be mapped such that only a subset of the nodes is used for the mapping. To do so, we use a recursive k-

²<https://www.open-mpi.org/projects/netloc>

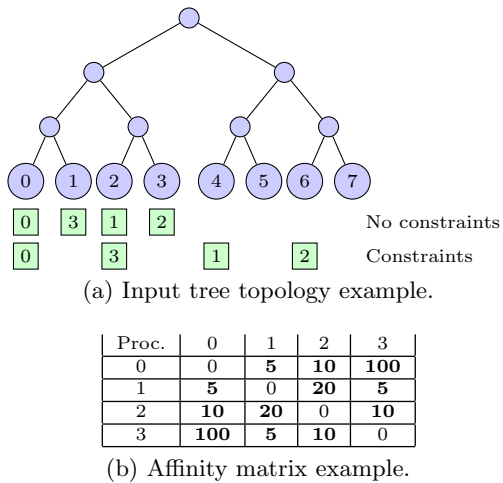


Figure 2: Example of TreeMatch output (green square) based on the affinity matrix and the tree topology. The first line is without constraints: in this case the hop-byte metric is 360. The second line is when only cores with even numbers are allowed to execute processes (hop-byte is 660 in this case)

partitioning algorithm where we add dummy processes that are forced to be mapped onto unavailable resources while real processes are mapped to actual available resources.

In Fig. 2, we describe an example where we map 4 processes on an architecture featuring 8 computing resources and structured as a 3-levels tree. We display 2 cases: one without constraints and the other where only cores with even numbers are available for mapping.

3.2 Job Allocation Strategy

We implement a new selection option for the SLURM `cons_res` plugin. In this case the regular best-fit algorithm used for nodes selection is replaced by TREEMATCH.

To this end we need to provide three pieces of information: a job affinity matrix, the cluster topology and the constraints due to other jobs allocations.

The communication matrix is provided at job submission time through a distribution option available in the `srun` command:

```
srun -m TREEMATCH=/comm/matrix/path cmd
#SBATCH -m TREEMATCH=/comm/matrix/path.
```

Its location (path) is then stored by the SLURM controller in the data structure describing a job and can be used by TREEMATCH for allocation.

As for the global cluster topology, it is provided to the controller by a new parameter in the configuration file: `TreematchTopologyFile=/topology/file/path`.

Whenever a job allocation is computed, this topology is completed by constraints informations. These constraints are provided by the nodes and cores bitmaps used by the SLURM controller to describe the cluster utilization. We need to translate this topology description into the TREEMATCH topology.

TREEMATCH considers computing units as selection granularity and assign them an id considering the global topology. It must be the same for the SLURM selection plugin using TREEMATCH. Hence we use the `cons_res` plu-

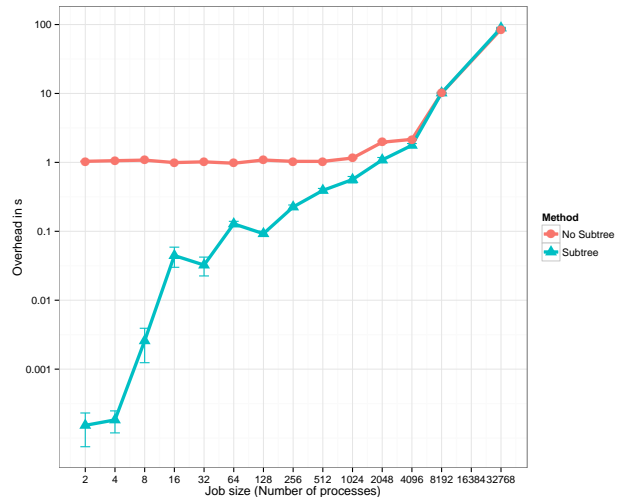


Figure 3: Comparison of TREEMATCH overhead (s) for different job size on a cluster with 80640 cores between methods with and without subtrees

gin with the configuration `SelectTypeParameters=CR_CPU` or `CR_Core`. In this case SLURM uses a cores bitmap describing precisely the location of unused CPUs inside nodes relatively to the nodes bitmap. Therefore, we need to translate SLURM local CPU ids into global TREEMATCH CPU ids. Then, we use the constraints feature of TREEMATCH (described in Section 3.1) to only use CPUs not already allocated to a running job. The CPUs chosen by TREEMATCH must then be translated again in new bitmaps for SLURM to use.

However, in the case of a large topology, our algorithm overhead increases: the larger the topology, the longer the TREEMATCH algorithm takes. To reduce this time, we also implement an alternative method which first finds a subtree in the global topology. Then, TREEMATCH uses this subtree to rapidly choose the job allocation. To find this subtree we search through the topology tree from the leaves up to the root and from left to right. We stop as soon as we find a node with enough unused CPUs. For instance, if we consider Fig. 1 and we assume that node `n0` is occupied instead of `n3`, then the first tree with 2 CPUs is `n1` and if we need 6 CPUs, we shall select subtree `t1`.

Fig. 3 compares the overhead of this algorithm with and without subtree utilization on a cluster featuring 80640 cores. It shows that, for jobs using less than 4096 cores, the subtree technique reduces the overhead. In any case both approaches takes less than 1s. At some point, the time increases linearly with the application size. However, as shown in the experiments (Sec. 4), the TREEMATCH overhead is largely compensated by the execution time gain. Moreover, for large applications, it is possible to compute the mapping at the node level (instead of computing it at the core level): hence a full-size application (80640 cores) requires 5040 nodes which leads to an overhead of a few seconds.

For the experiments described in Section 4 we need to modify the jobs run times dynamically according to their allocation. To do this we compute for each job both the SLURM allocation and the TREEMATCH one. Then we compute R , the ratio between their hop-byte cost (c.f. Section 3.1). We

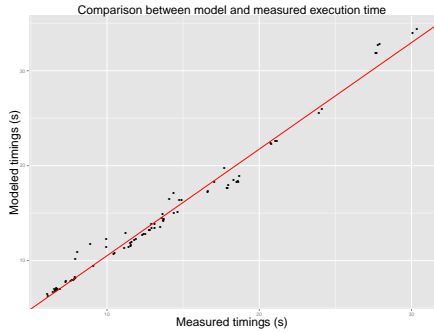


Figure 4: TREEMATCH measured time vs. modeled time for the minighost application with a communication ratio between 5% and 45%

model job runtimes with computation times and communication times: $T = T_{calc} + T_{comm}$. Let α be the ratio of communication time of the whole runtime: $T_{comm} = \alpha T$. Hence, $T = \alpha T + (1 - \alpha)T$. TREEMATCH impacts only the communication cost. Therefore, we model the execution time T' using the TREEMATCH allocation with:

$$\begin{aligned} T' &= T_{calc} + RT_{comm} = R\alpha T + (1 - \alpha)T \\ &= (1 + R\alpha - \alpha)T \end{aligned}$$

We validate this model with the minighost application [3] that computes a stencil in various dimensions. We execute 84 runs with various settings (number of processors, different parameters) using a round-robin placement or a mapping computed with TREEMATCH. The minighost output also provides the percentage of communication in a run. In our case, this ranges from 5% to 45%. Fig. 4 shows the validation of the above model. On the x-axis is the TREEMATCH runtime and on the y-axis is the predicted time based on the ratio R of the hop-byte of the TREEMATCH mapping and the SLURM mapping, α the percentage of communication and T the measured execution runtime. We see a very strong correlation between both timings even though the modeled timings tend to be slightly larger than the real ones.

4. EXPERIMENTAL VALIDATION

4.1 Emulation Experimental Setup

Our experiments have been carried out on the Edel cluster from the Grid'5000 Grenoble site. Edel is composed of 72 nodes featuring 2 Intel Xeon E5520 CPUs (2.27 GHz, 4 cores/CPU) and 24GB of memory.

We emulate Curie (a TGCC cluster with 5040 nodes and 80640 cores³) using a SLURM internal emulation technique called `multiple-slurmd` initially described and used in [10]. SLURM uses daemons: one `slurmctld` as the controller and one `slurmd` on each node. To emulate a larger cluster, we use 16 Edel nodes and launch 315 `slurmd` daemons on each node. We can consequently submit jobs as if we were working on the Curie cluster, emulating all the job scheduling overheads. We use simple jobs (just performing a call to `sleep`) in order to provide the necessary time and space illusion to the controller that a real job is actually executing.

³<http://www-hpc.cea.fr/en/complexe/tgcc-curie.htm>

We base our experiments on a Curie workload trace taken from the Parallel Workload Archive⁴. We have two sets of jobs. The first one is to fill the cluster, and the jobs belonging to this set are always scheduled using SLURM in order to have the same starting point for all the experiments. The second set, called the *workload*, is the one we actually use to compare the different strategies.

All the measurements are done through the SLURM logging system which gives us workload traces similar to the ones we obtain from Curie.

Finally, to use TREEMATCH we need to provide each job with a communication matrix. For these experiments we use randomly generated matrices featuring various sparsity rates. Indeed, the name of the application does not appear in the workload trace and therefore we cannot know the gain an optimized mapping would yield for a given entry in the trace. However, this depends on the communication ratio of the application (the higher this ratio, the larger is the possible gain in terms of runtime) and its communication pattern. Here, based on recent results [8, 11], we design the matrices such that the gain is similar to real-world applications. On average, the observed gain is 4% (resp. 11% and 18%) for a communication ratio of 10% (resp. 30% and 50%).

4.2 Emulation Results

We compare 4 cases : the classical topology-aware SLURM selection (SLURM), the same but using TREEMATCH for process placement after the allocation process and just before the execution starts (TM-A), TREEMATCH used both for the allocation process and for the process placement (TM-I) and finally the same but using the subtree technique to reduce the overhead (TM-Isub).

To evaluate our results, we use several metrics (two are for the whole workload and two are for each individual job):

- makespan: this is the time taken between the submission of the first job and the completion of the last job of the *workload*.
- utilization: this is the ratio between the CPUs used and the total number of CPUs in the cluster during the execution of the *workload*.
- job flowtime (or turnaround time): this is the time between the submission and the completion of a given job.
- job runtime: this is the time between the start and the completion of a given job.

In our case, the *workload* comprises 60 jobs. To keep the duration reasonable we decrease the jobs runtimes by a 50% factor. Figure 5 describes the results obtained for this workload and two values of α (1/3 and 1/2). Figure 5a shows that using TREEMATCH to reorder the process ranks reduces the makespan but using it inside SLURM to allocate nodes decreases it even more. This is what is shown in Fig. 1: incorporating TREEMATCH in SLURM gives more room for optimization as the mapping is not constraints by the allocation. Moreover, the subtree optimization leads to comparable results than without the optimization. This is due to the fact that in this case the makespan is determined by a small set of jobs and hence the impact of this optimization

⁴<http://www.cs.huji.ac.il/labs/parallel/workload/>

Com	SLURM	TM-A	TM-Sub	TM-I
50%	8318	6407	6073	6077
33%	8316	7502	6821	6887

(a) Makespan

Com	SLURM	TM-A	TM-Sub	TM-I
50%	33%	42%	44%	44%
33%	33%	36%	40%	39%

(b) Utilization

Figure 5: Workload Metrics for the different strategies and different amount of communication ratio

33% of communication

TM-Sub	22.50 s / 1.16	38.20 s / 1.44	205.85 s / 1.09
[5 s, 18 s]	TM-I	15.70 s / 1.24	183.35 s / 0.95
[14 s, 20 s]	[7 s, 13 s]	TM-A	167.65 s / 0.76
[20 s, 253 s]	[6 s, 213 s]	[4 s, 185 s]	SLURM

(a) 33% of communication

33% of communication

TM-Sub	10.20 s / 1.19	47.83 s / 1.47	322.23 s / 1.27
[4 s, 14 s]	TM-I	37.63 s / 1.24	312.03 s / 1.06
[12 s, 23 s]	[3 s, 11 s]	TM-A	274.40 s / 0.86
[27 s, 396 s]	[13 s, 306 s]	[11 s, 307 s]	SLURM

(b) 50% of communication

Figure 6: Statistical comparison of selection methods: flow time

is not visible for this metric. We also see that the larger the communication ratio the greater the gain, this is expected as TREEMATCH optimizes communication only. This is tested through simulation in Section 4.3.

Figure 5b also shows that for the same submission workload, TREEMATCH improved the resource utilization.

In Fig. 6 and Fig. 7, we use paired comparisons between different strategies for respectively jobs flowtime and jobs runtime. Here, we consider job-wise metrics, therefore we want to understand if, when we average all the jobs, a strategy turns out to be better than another. Each strategy is displayed on the diagonal. On the upper right, we have the average difference between the strategy on the column and the one on the row and the geometric mean of the ratios. For instance, in Fig. 6a, we see that on average the job flowtime is 183.35s faster with TM-I than with SLURM and the average ratio is 0.95. On the lower left part, we plot the 90% confidence interval of the corresponding mean. The interpretation is the following: if the interval is positive, then the strategy on the row is better than the strategy on the line with a 90% confidence. In this case, the correspond-

33% of communication

TM-Sub	18.90 s / 1.20	37.45 s / 1.54	200.17 s / 1.08
[4 s, 14 s]	TM-I	18.55 s / 1.28	181.27 s / 0.90
[14 s, 20 s]	[9 s, 14 s]	TM-A	162.72 s / 0.70
[13 s, 252 s]	[3 s, 212 s]	[-2 s, 176 s]	SLURM

(a) 33% of communication

33% of communication

TM-Sub	7.03 s / 1.22	48.43 s / 1.54	317.10 s / 1.17
[3 s, 11 s]	TM-I	41.40 s / 1.27	310.07 s / 0.96
[13 s, 23 s]	[6 s, 13 s]	TM-A	268.67 s / 0.76
[22 s, 383 s]	[6 s, 303 s]	[2 s, 305 s]	SLURM

(b) 50% of communication

Figure 7: Statistical comparison of selection methods: runtime

ing mean is highlighted in green. If the interval is negative the strategy on the line is better than the one on row and the corresponding mean is highlighted in red. Otherwise, we cannot statistically conclude with a 90% confidence on which strategy is the best and we do not highlight the corresponding mean. For example, on Figure 6a we can see that using TREEMATCH in SLURM is better than not using it. Moreover, here we see that using the subtree optimization improves the metric. For all the cases we see that TM-Sub is better than TM-I that is better than TM-A. Therefore, restricting the usage of TREEMATCH improves the performance as the gain in computing a solution overcome the loss in terms of quality of this solution.

Moreover, both flowtime and runtime using TREEMATCH in SLURM are shorter than using TREEMATCH after SLURM, with a ratio between 1.44 and 1.54. We can also see that the more an application communicates, the smaller are the average gaps. For example, between TM-I and TM-Sub (with a 33% of communication ratio), the average difference is 22.5 s, but for a 50% ratio it is 10.2s. In these experiments, the cluster is already full when submitting the first jobs. Therefore, a part of their flowtime corresponds to the wait for a free allocation.

Figure 7 shows the comparison of jobs runtimes. We observe similar behavior except that the confidence interval between SLURM and TM-A does not allow to conclude with 90% confidence that TM-A is better than SLURM.

Through these experiments we observe that using TREEMATCH in the allocation process induces no negative effects and improves the global use of a cluster. Moreover, from a user point of view, using TREEMATCH can also be profitable by decreasing the runtime of his/her jobs.

4.3 Simulation Results

As the experiments done in the above section are carried

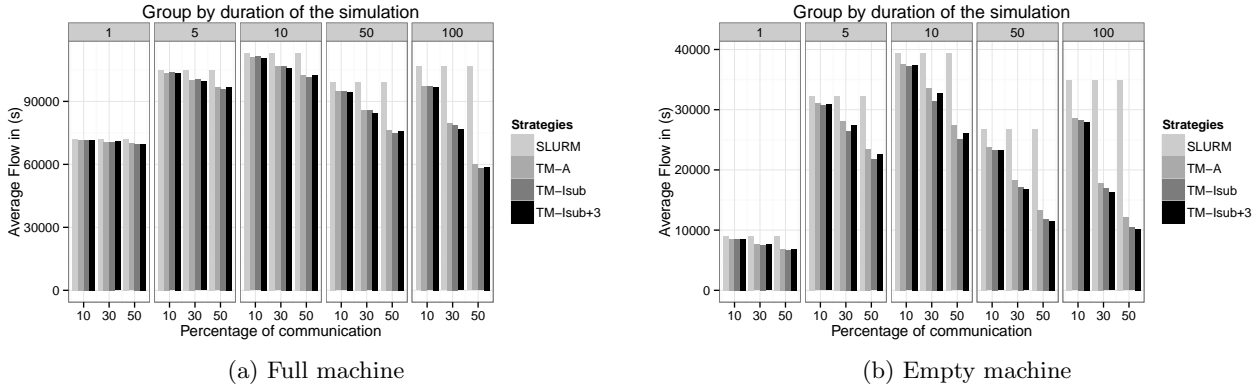


Figure 8: Average flowtime using the Curie trace with different strategies and various percentage of communication

out through emulation, they are very long to compute (as long as the real execution times). In order to cover a larger set of test cases and a longer time-scale we design a simulator that simulates both the job selection part and the job execution part. Our simulator is an event-based one that reads the machine topology and the job submission trace workload and computes the start time and end time of each job based on its duration (given by the workload) and the allocation. For the job selection step, we implement the same algorithm than we use in the above section and the time to compute the allocation is based on the duration of TREEMATCH when it is used or is set to 2 seconds when SLURM is used. For the execution time part, we use the formula shown in Section 4.1 if we use TREEMATCH. Otherwise, the duration given by the Curie trace is used. It is accurate enough to provide makespan duration with an average absolute error of less than 3% for Figure 5a.

The experiments described here represent more than 12 millions node-hours. Note that allocations above 1 million node-hours are only given to application projects proposal by supercomputing centers and never to software stack optimization projects⁵, hence justifying the use of simulation.

Figure 8a shows the average flow of the jobs in the case where the machine is already full of jobs (that have all been scheduled using the regular SLURM strategy). We group the measurements by simulation duration (i.e. for group 50, we consider only the jobs submitted during the first 50 hours): we go from 1 hour (365 jobs) to 100 hours (13687 jobs). On the x-axis we display the different percentage of communication (from 10% to 50%) and we have 4 strategies: the plain SLURM, TREEMATCH applied at the beginning of the job to map processes to resources after SLURM has allocated the nodes (TM-A), and TREEMATCH used in SLURM to compute the allocation and the mapping using the minimal subtree (TM-Isub), or 3 levels above the minimal subtree (TM-Isub+3). We see that the impact of TREEMATCH on the flow increases with the duration of the simulation because at the beginning of the simulation, as the machine is full, the flowtime depends mainly on the time a job has to wait before starting, while as time goes the impact on the improvement of the mapping due to TREEMATCH accumulates. Here, we do not see much difference between the

different strategies involving TREEMATCH because there is less room for optimization when the machine is fully utilized. However, the results are consistent, in terms of quality with the emulation results presented in the previous section.

Figure 8b shows the average flow of the jobs in the case where the machine is totally empty. In this case, we see that the gain with TREEMATCH increases and appears earlier which corroborates the hypothesis made in the previous paragraph. Moreover, as we have more opportunities for optimization we see that using TREEMATCH in SLURM is more beneficial than using it just before the job execution. We also see a large gap for hour 5 because in the workload a large job (32768 cores) is submitted at 8461s, that takes a long time to schedule and that uses a substantial part of the machine, thus impacting all the subsequent jobs.

As in the previous section, we see that even with a small average gain on each jobs (4% for the 10% communication ratio case to 17% for the 50% communication ratio), we are able to achieve very large gain on the overall. This is due to the fact that these gains accumulates during the workload lifetime. We therefore expect even greater gains on real settings as the operational lifetime of a real machine is much longer than the experiments done here.

5. RELATED WORKS AND DISCUSSION

The idea of using the most adequate hardware resource to a specific application is not new and has been explored in previous work. It has been particularly popular in the context of grids environments ([17], [27], [25]) where it is important to select the best set of resources (clusters in this case) to use. Such work try to reduce the impact of WAN communication in grids but do not address the deeper details of the physical topology, such as NUMA effects or cache hierarchy for instance.

More recently, some works have targeted a specific type of applications, that is, MapReduce-based applications. For instance, the TARA [16] uses a description of the application to allocate the resources. However, this work is tailored for a very specific class of applications and does not address hardware details.

The mapping of a parallel applications' tasks to the physical processors based on the network topology can lead to important performance improvements [4]. Network topology characteristics can be taken into account by the scheduler [20] so as to favor the choice of group of nodes that are

⁵See PRACE core hours award for instance: <http://www.prace-ri.eu/hermit-awardees>

placed on the same network level, connected under the same network switch or even placed close to each other so as to avoid long distance communications. This kind of feature is taken into account by most of open-source and proprietary RJMSs. However even if most of them use the characteristics of the underlying physical topology, they eventually fail to take into consideration the application behaviour when allocating resources and this is something that this work specifically addresses. HTCondor (formerly Condor) leverages a so-called *matchmaking* approach [23] that allows it to match the applications needs to the available hardware resources. However, the application *behaviour* is not part of this match-making and HTCondor targets both clusters and networks of workstations. SLURM [30], as previously described, provides an option to minimize the number of network switches used in the allocation, so as to reduce the communication costs during the application execution (switches that are the deeper in the tree topology are supposed to be the less costly than upper ones). The same idea of topology-aware placement is exploited by PBS Pro [22], Grid Engine[21], and LSF [26]. Fujitsu [9] provides the same but only for its proprietary Tofu network. As far as our knowledge, SLURM [30] remains the only one providing a *best-fit* topology-aware selection whereas the others propose *first-fit* algorithms.

Some other RJMS offer task placement options that can enforce a clever placement of the application processes. That is the case of Torque [7] which proposes a NUMA-aware job task placement. OAR [6] uses a flexible hierarchical representation of resources which offers the possibility to place the application processes upon the hierarchy within the computing node. However, in these existing works, only the network topology is taken in account and the nodes internal architecture is left unaddressed when performance gains are expected from exploiting the memory hierarchy.

Jingjin Wu et al. in [28] introduced a hierarchical task mapping strategy for modern supercomputers based on generic recursive algorithms for both fat-tree and torus network topologies showing very good performance with low overhead. Rashti et al. [24] proposed a weighted graph model for the whole physical topology of the computing system, including both the inter and intra node topologies. Even if both previous related works have shown interesting results with application sets, they have not been integrated with real resource and job management system neither tested with real workload traces while souissions mixesch is our case in this paper.

A study for torus network topology [1] showed how processor ordering takes place based on space filling curve, such as Hilbert Curve, to map the nodes of the torus onto a 1-dimensional list in order to preserve locality information. This paper described the study about the allocation strategies implemented on the proprietary Cray Application Level Placement Scheduler (ALPS). Similar strategies, have been recently incorporated within SLURM with⁶ (or without⁷) the use of ALPS. Another interesting work [29] adapted only for torus topology, presented a window-based locality-aware job scheduling strategy that tries to optimize job and system performance in the same time. Its goal is to preserve node contiguity by considering multiple jobs for scheduling while making use of the 0-1 Multiple Knapsack problem for

resource allocation. The last 2 related works do not consider communication patterns as parameters within the algorithms.

Several binding policies are available, and they are compatible with the policies implemented in Open MPI. In all these solutions, the user has to retrieve the architectural details before submitting his/her job. Also, the placement options offered leave the user with the burden to determine his/her policy beforehand, and the application communication scheme is not taken into account.

In our case, we improve this functioning on three levels: first, we take into account not only the network but also the node internal structure. The information used is based on the *structure* of the nodes and the memory hierarchy. In other words, we do not use latency and bandwidth figures to compute our allocation. Then, this information is retrieved directly by our plugin does not have to be supplied by the user. All the technical details are hidden. Last, but not least, we also take into account not only the architecture but also the application behaviour both for the allocation and the execution of a job.

6. CONCLUSIONS

Job scheduling plays a crucial role in cluster administration, enabling both better response time and resource usage. In this paper, we tackle the problem of allocating and mapping jobs according to a cluster topology and application process affinity. We extend TREEMATCH to design a new allocation policy that allocates and maps at the same time application processes on the resources, based on the communication matrix of the considered application. Such strategy is implemented in the SLURM `cons_res` plugin. We test this strategy on emulation and simulation and compare it with the standard SLURM topology-aware policy and the method consisting in mapping processes after the allocation is determined.

Results show that taking into account application characteristics and the topology provides better makespan, flow time, utilization and job runtime compared to the standard topology-aware and compact SLURM policy. We also show that the level at which we consider the topology impacts the performance. It is better to have a more local view of the topology than only a global view since in this latter case, allocation quality is slightly better but longer to compute. Last, even if not all the jobs are able to use this strategy all of them benefit from it with a reduced flowtime.

For future work, we would like to investigate the following research axes. First, we would like to look at fragmentation metrics. Indeed, the way jobs are allocated impacts the global resource usage and this aspect should be quantified. Also, we would like to find means to gather in a systematic fashion applications communication patterns in order to create an applications classification based on these patterns and then implement this solution in production. We would also like to validate this approach in other job scheduler such as OAR [6]. Concerning SLURM integration and extensions, we are currently working on the inclusion of our new developments in the next official SLURM release.

7. ACKNOWLEDGMENTS

Experiments presented in this paper were carried out using the Grid'5000 testbed, supported by a scientific inter-

⁶http://slurm.schedmd.com/cray_alps.html

⁷<http://slurm.schedmd.com/cray.html>

est group hosted by Inria and including CNRS, RENATER and several Universities as well as other organizations (see <https://www.grid5000.fr>). Part of this work is also supported by the ANR MOEBUS project ANR-13-INFR-0001. This work is partially funded under the ITEA3 COLOC project #13024.

8. REFERENCES

- [1] C. Albing, N. Troullier, S. Whalen, R. Olson, J. Glenski, H. Pritchard, and H. Mills. Scalable node allocation for improved performance in regular and anisotropic 3d torus supercomputers. In *Recent Advances in the Message Passing Interface - 18th European MPI Users' Group Meeting, EuroMPI 2011, Santorini, Greece, September 18-21, 2011. Proceedings*, pages 61–70, 2011.
- [2] S. M. Balle and D. J. Palermo. Enhancing an open source resource manager with multi-core/multi-threaded support. In *Job Scheduling Strategies for Parallel Processing, 13th International Workshop, JSSPP 2007, Seattle, WA, USA, June 17, 2007. Revised Papers*, pages 37–50, 2007.
- [3] R. F. Barrett, C. T. Vaughan, and M. A. Heroux. Minighost: a miniapp for exploring boundary exchange strategies using stencil computations in scientific parallel computing. *Sandia National Laboratories, Tech. Rep. SAND2011-5294832*, 2011.
- [4] A. Bhatele, E. J. Bohm, and L. V. Kalé. Topology aware task mapping techniques: an api and case study. In *PPOPP*, pages 301–302, 2009.
- [5] F. Broquedis, J. Clet-Ortega, S. Moreaud, N. Furmento, B. Goglin, G. Mercier, S. Thibault, and R. Namyst. Hwloc: a Generic Framework for Managing Hardware Affinities in HPC Applications. In *Proceedings of the 18th Euromicro International Conference on Parallel, Distributed and Network-Based Processing (PDP2010)*, Pisa, Italia, Feb. 2010. IEEE Computer Society Press.
- [6] N. Capit, G. Da Costa, Y. Georgiou, G. Huard, C. Martin, G. Mounié, P. Neyron, and O. Richard. A batch scheduler with high level components. In *Cluster computing and Grid 2005 (CCGrid05)*, Cardiff, United Kingdom, 2005. IEEE.
- [7] A. computing. Torque resource manager. <http://docs.adaptivecomputing.com/torque/6-0-0/Content/topics/torque/2-jobs/monitoringJobs.htm>.
- [8] E. H. Cruz, M. Diener, L. L. Pilla, and P. O. Navaux. An efficient algorithm for communication-based task mapping. In *Parallel, Distributed and Network-Based Processing (PDP), 2015 23rd Euromicro International Conference on*, pages 207–214. IEEE, 2015.
- [9] Fujitsu. Interconnect topology-aware resource assignment. <http://www.fujitsu.com/global/Images/technical-computing-suite-bp-scl2.pdf>.
- [10] Y. Georgiou and M. Hautreux. Evaluating scalability and efficiency of the resource and job management system on large HPC clusters. In *Job Scheduling Strategies for Parallel Processing, 16th International Workshop, JSSPP 2012, Shanghai, China, May 25, 2012. Revised Selected Papers*, pages 134–156, 2012.
- [11] T. Hoefler and M. Snir. Generic Topology Mapping Strategies for Large-Scale Parallel Architectures. In *ICS*, pages 75–84, 2011.
- [12] E. Jeannot, E. Meneses, G. Mercier, F. Tessier, and G. Zheng. Communication and topology-aware load balancing in charm++ with treematch. In *IEEE Cluster*, page 8, Indianapolis, IN, USA, Sept. 2013.
- [13] E. Jeannot and G. Mercier. Near-Optimal Placement of MPI processes on Hierarchical NUMA Architectures. In Pasqua D’Ambra, Mario Rosario Guarracino, and Domenico Talia, editors, *Euro-Par 2010 - Parallel Processing, 16th International Euro-Par Conference*, volume 6272 of *Lecture Notes on Computer Science*, pages 199–210, Ischia Italie, SEPT 2010. Springer.
- [14] E. Jeannot and G. Mercier. Near-optimal placement of mpi processes on hierarchical numa architectures. *Euro-Par 2010-Parallel Processing*, pages 199–210, 2010.
- [15] E. Jeannot, G. Mercier, and F. Tessier. Process Placement in Multicore Clusters: Algorithmic Issues and Practical Techniques. *IEEE Trans. Parallel Distrib. Syst.*, 25(4):993–1002, 2014.
- [16] G. Lee, N. Tolia, P. Ranganathan, and R. H. Katz. Topology-aware resource allocation for data-intensive workloads. In *Proceedings of the First ACM Asia-pacific Workshop on Workshop on Systems, APSys '10*, pages 1–6, New York, NY, USA, 2010. ACM.
- [17] C. Liu, L. Yang, I. Foster, and D. Angulo. Design and evaluation of a resource selection framework for grid applications. In *Proceedings of the 11th IEEE International Symposium on High Performance Distributed Computing, HPDC '02*, pages 63–, Washington, DC, USA, 2002. IEEE Computer Society.
- [18] G. Mercier and J. Clet-Ortega. Towards an Efficient Process Placement Policy for MPI Applications in Multicore Environments. In *EuroPVM/MPI*, volume 5759 of *Lecture Notes in Computer Science*, pages 104–115, Espoo, Finland, Sept. 2009. Springer.
- [19] G. Mercier and E. Jeannot. Improving MPI Applications Performance on Multicore Clusters with Rank Reordering. In *EuroMPI*, volume 6960 of *Lecture Notes in Computer Science*, pages 39–49, Santorini, Greece, Sept. 2011. Springer.
- [20] J. Navaridas, J. Miguel-Alonso, F. J. Ridruejo, and W. Denzel. Reducing complexity in tree-like computer interconnection networks. *Parallel Computing*, 36(2-3):71–85, 2010.
- [21] Oracle. Grid engine.
- [22] PBSWorks. Pbs. <http://www.pbsworks.com/PBSProduct.aspx?n=PBS-Professional&c=Overview-and-Capabilities>.
- [23] R. Raman, M. Livny, and M. Solomon. Matchmaking: Distributed resource management for high throughput computing. In *Proceedings of the Seventh IEEE International Symposium on High Performance Distributed Computing (HPDC7)*, Chicago, IL, July 1998.
- [24] M. J. Rashti, J. Green, P. Balaji, A. Afsahi, and W. Gropp. Multi-core and network aware MPI topology functions. In *Recent Advances in the Message Passing Interface - 18th European MPI Users' Group*

- Meeting, EuroMPI 2011, Santorini, Greece, September 18-21, 2011. Proceedings*, pages 50–60, 2011.
- [25] C. A. Santos, A. Sahai, X. Zhu, D. Beyer, V. Machiraju, and S. Singhal. *Utility Computing: 15th IFIP/IEEE International Workshop on Distributed Systems: Operations and Management, DSOM 2004, Davis, CA, USA, November 15-17, 2004. Proceedings*, chapter Policy-Based Resource Assignment in Utility Computing Environments, pages 100–111. Springer Berlin Heidelberg, Berlin, Heidelberg, 2004.
- [26] C. Smith, B. McMillan, and I. Lumb. Topology aware scheduling in the lsf distributed resource manager. In *Proceedings of the Cray User Group Meeting*, 2001.
- [27] O. Sonmez, H. Mohamed, and D. Epema. Communication-aware job placement policies for the koala grid scheduler. In *Proc. of the Second IEEE International Conference on e-Science and Grid Computing (e-Science'06)*, pages 79–86. IEEE Computer Science, Dec 2006.
- [28] J. Wu, X. Xiong, and Z. Lan. Hierarchical task mapping for parallel applications on supercomputers. *The Journal of Supercomputing*, 71(5):1776–1802, 2015.
- [29] X. Yang, Z. Zhou, W. Tang, X. Zheng, J. Wang, and Z. Lan. Balancing job performance with system performance via locality-aware scheduling on torus-connected systems. In *2014 IEEE International Conference on Cluster Computing, CLUSTER 2014, Madrid, Spain, September 22-26, 2014*, pages 140–148, 2014.
- [30] A. Yoo, M. Jette, and M. Grondona. Slurm: Simple linux utility for resource management. In D. Fietelson, L. Rudolph, and U. Schwiegelshohn, editors, *Job Scheduling Strategies for Parallel Processing*, volume 2862 of *Lecture Notes in Computer Science*, pages 44–60. Springer Berlin Heidelberg, 2003.
- [31] H. Yu, I.-H. Chung, and J. Moreira. Topology mapping for blue gene/l supercomputer. In *Proceedings of the 2006 ACM/IEEE Conference on Supercomputing, SC '06*, New York, NY, USA, 2006. ACM.