

A Set Coverage-based Mapping Heuristic for Scheduling Distributed Data-Intensive Applications on Global Grids

Srikumar Venugopal and Rajkumar Buyya

Grid Computing and Distributed Systems (GRIDS) Laboratory
Department of Computer Science and Software Engineering
The University of Melbourne, Australia
{srikumar, raj}@csse.unimelb.edu.au

Abstract—Data-intensive Grid applications need access to large datasets that may each be replicated on different resources. Minimizing the overhead of transferring these datasets to the resources where the applications are executed requires that appropriate computational and data resources be selected. In this paper, we introduce a heuristic for the selection of resources based on a solution to the Set Covering Problem (SCP). We then pair this mapping heuristic with the well-known MinMin scheduling algorithm and conduct performance evaluation through extensive simulations.

I. INTRODUCTION

Grids [1] aggregate computational, storage and network resources to provide pervasive access to their combined capabilities. In addition, Data Grids [2], [3] provide services such as low latency transport protocols and data replication mechanisms to distributed data-intensive applications that need to access, process and transfer large datasets stored in distributed repositories. Such applications are commonly used by communities of researchers in domains such as high-energy physics, astronomy and biology.

The work in this paper is concerned with scheduling data-intensive applications that can be considered as a collection of tasks without interdependencies, each of which requires multiple datasets, onto a set of Grid resources. An astronomy image-processing application following this model is described by Yamamoto, et al. [4]. Each task is translated into a job that is scheduled on to a computational resource and requests datasets from the storage resources (or *datahosts*). Each of these datasets may be replicated at several locations that are connected to each other and to the computational sites (or *compute resources*) through networks of varying capability. This scenario is illustrated in Figure 1. For a job, the scheduling strategy has to select a set of resources consisting of a compute resource for execution and a subset of data hosts such that each dataset required for the job can be accessed from one of the data hosts in the set. We term this as the *mapping* problem and in this paper, present a heuristic based on a solution to the SCP. We evaluate it against other heuristics through extensive simulations.

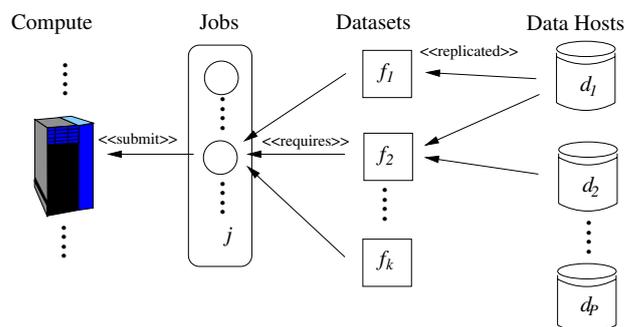


Fig. 1: Mapping Problem.

The rest of the paper is structured as follows: we present the related work followed by the resource model and the application model that we target in the research presented in this paper. The mapping heuristic is presented in the following section and is succeeded by details of experimental evaluation and the consequent results. Finally, we present our conclusions.

II. RELATED WORK

Previous publications ([5]–[8]) on scheduling in Data Grids have been more concerned with the relationship between job assignment and data replication. That is, they assign a single data-intensive job to a compute resource based on how close it is to the data required and if the data is not available at that resource, then initiate replication. Casanova, et al [9] introduce a new heuristic *XSufferage* – based on the *Sufferage* heuristic in [10] – that takes into consideration sharing of files between tasks. For a task, it favors a resource where the data is already present over one where it has to be copied. However, the source of all the files for the tasks is the resource that dispatches the jobs. Mohamed and Epema [11] present a Close-to-Files algorithm for a similar application model, though restricted to one dataset per job, that searches the entire solution space for a combination of computational and storage resources to minimize execution time.

The application model in this paper is, however, closer to that of Giersch, et. al [12] who consider the general problem of scheduling tasks requiring multiple files that are available from multiple sources and prove that this problem is NP-complete. However, the approach followed in their paper is that of scheduling the jobs first and then replicating the data so as to minimize access time. Khanna, et. al [13] propose a hypergraph-based approach for scheduling a set of similar tasks with a view to minimize the I/O overhead by considering the sharing of files between the tasks. However, they do not take into account the aspect of replication as the files have only single sources.

The work in this paper is distinct from the related work because it considers: a) the problem of selecting a resource set for a job requiring multiple datasets in an environment where the data is available from multiple sources due to prior replication and b) the selection of computational and data resources for such a resource set to be interconnected. In a previous publication [14], we have introduced a greedy mapping heuristic for deadline and budget-based scheduling of applications following this model. In this paper however, we concentrate on the objective of reducing the total makespan of such applications.

III. MODEL

A. Resource Model

We model the target data-intensive computing environment based on existing production testbeds such as the European DataGrid testbed [3] or the United States Grid3 testbed [15]. As an example, Figure 2 shows a subset of European DataGrid Testbed 1 derived from Bell, et. al [6]. The resources in the figure are spread across 7 countries and belong to different autonomous administrative domains.

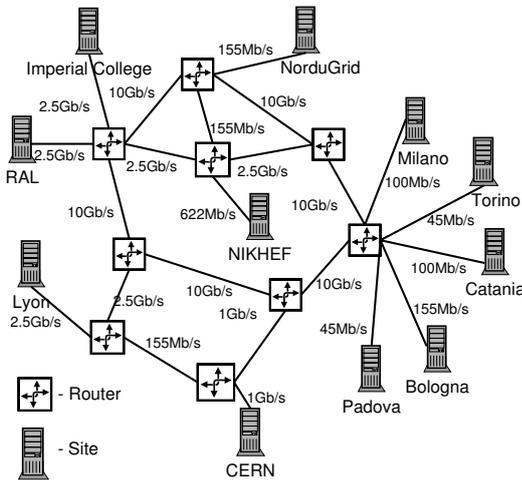


Fig. 2: European Data Grid Testbed 1 [6].

In such Grid networks, we consider a data-intensive computing environment to consist of a set of M compute resources $R = \{r_m\}_{m=1}^M$ and a set of P data hosts, $D = \{d_p\}_{p=1}^P$. A

compute resource is a high performance computing platform such as a cluster consisting of processing nodes that are connected in a private local area network and are managed by a batch job submission system hosted at the “head” node that is connected to the public Internet. A data host can be a storage resource such as a Mass Storage Facility connected to the Internet or may be simply a storage device attached to a compute resource in which case, it inherits the network properties of the latter. It is important to note that even in the second case, the data host is considered as a separate entity from the compute resource. We consider the logical network topology wherein each resource is connected to every other resource by a distinct logical network link whose capacity is given by the bottleneck bandwidth of the physical network between the resources. The time taken by a compute resource to access a dataset located on the storage resource at the same site is limited only by the intra-site bandwidth if the storage is a separate physical machine or by the bandwidth between the hard disk and other peripherals if the storage is on the compute machine itself. In both cases, it is considered to be an order of magnitude lower than the time taken to access a dataset through the Internet from other sites as there is contention for bandwidth among the various sites.

Data is organised in the form of datasets that are replicated on the data hosts by a separate replication process that follows a strategy (e.g. [6]) that takes into consideration various factors such as locality of access, load on the data host and available storage space. Information about the datasets and their location is available through a catalog such as the Storage Resource Broker Metadata Catalog [16].

B. Application Model

The application is composed of a set of N jobs without interdependencies, $J = \{j_i\}_{i=1}^N$. Typically, $N \gg M$. Each job $j, j \in J$ requires a subset $F^j = \{f_k^j\}_{k=1}^K$ of a set of datasets, F , which are each replicated on a subset of P data hosts, $D = \{d_p\}_{p=1}^P$. For a dataset $f \in F$, $D_f \subseteq D$ is the set of data hosts on which f is replicated. These sets of data hosts are not pairwise disjoint, that is, a data host can provide any number of datasets at the same time. For a job j , a *resource set* is denoted by $S^j = \{\{r\}, \{d_l\}_{l=1}^L\}$ where $r \in R$ is the compute resource selected for executing the job and $d_l \in D_{f_k^j}$ is the data host selected for accessing $f_k^j \in F^j$. Since multiple datasets can be retrieved from one data host, $L \leq K$, the number of datasets required by the job.

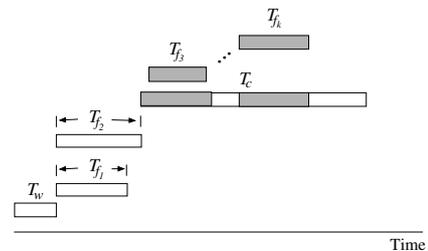


Fig. 3: Job Execution Stages and Times.

Figure 3 shows a generic example of a data-intensive job with the times involved in various stages shown along a horizontal time-axis. The gray colors show overlaps between computation and data transfer operations. T_w is the time spent in waiting in the queue on the compute resource and T_c is the time spent by the job in purely computational operations (also called computation time). T_w and T_c are functions of the load and processing speed of the compute resource. T_{f_i} is the time required to transfer the file f_i from its data host to the compute resource and is dependent on the available bandwidth between the two. The completion time for the job, T_j , is the wallclock time taken for the job to finish execution and is a function of T_w , T_c and T_{f_i} . For large datasets, the data transfer time impacts the completion time significantly. While the transfer time is determined by the manner in which the dataset is processed by the job, it is also influenced by the selection of data hosts. For example, many applications request and receive required datasets in parallel before starting computation. In this case, $T_j = T_w + \max_{1 \leq i \leq K}(T_{f_i}) + T_c$. However, the number of simultaneous transfers determines the bandwidth available at the receiver end for each transfer and therefore, the T_{f_i} . Transfer times can be minimized by locating a compute resource associated with a data host that has the maximum number of datasets required by the job so that the bulk of the data access is local. This would also benefit the case where the job accesses datasets sequentially.

Our aim here is to select a resource set that produces the Minimum Completion Time (MCT) for a job. We adopt the strategy of finding the resource set with the least number of data hosts required to access the datasets required for a job and then, finding a suitable compute resource to execute it. In doing so, we try to maximise the local access of datasets and thus, reduce the transfer times. We experimentally show that this approach produces schedules that are competitive with the best and is reasonably fast as well.

IV. SCHEDULING

The scheduler forms a part of a larger application execution framework such as a resource broker (e.g. [17], [18]). The resource broker is able to identify resources that meet minimum requirements of the application such as architecture (instruction set), operating system, storage threshold and data access permissions and these are provided as suitable candidates for job execution to the scheduler.

Figure 4 lists a generic scheduling algorithm for scheduling a set of jobs on a set of distributed compute resources. Each of the steps can be implemented independently of each other and therefore, many strategies are possible. In this paper, we concentrate on the process within the *for* loop, that is, finding the appropriate resource set for a job.

A. Mapping Heuristic

For a job $j \in J$, consider a graph $G^j = (V, E)$ where $V = (\bigcup_{f \in F_j} \{D_f\}) \cup F^j$ and E is the set of all directed edges $\{d, f\}$ such that $d \in D_f$. Our aim here is to find the minimum set H of data hosts such that there exists an edge

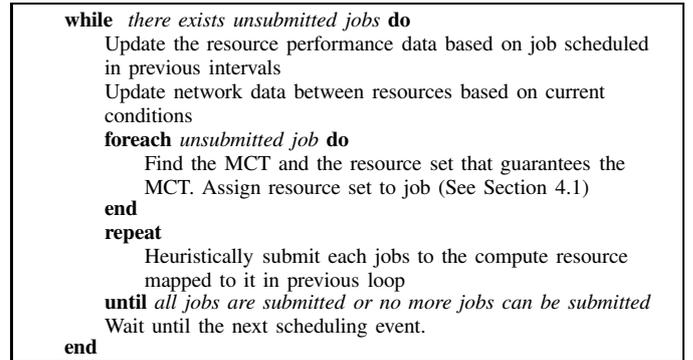


Fig. 4: A Generic Scheduling Algorithm.

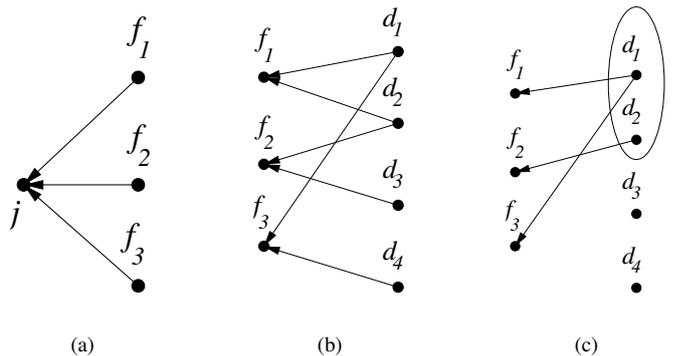


Fig. 5: (a) Job j dependent on 3 datasets. (b) Directed graph of data resources and data sets for job j . (c) A dominating set for the data graph.

from a member of H to f for every $f \in F_j$ in G and no other set $H' \subset H$ satisfies that requirement. The set H is called the *minimal dominating vertex set* for graph G . However, it is possible that more than one such set exists for a graph. Our interest is in finding a minimal dominating set of datahosts such that the completion time for j is reduced. Figures 5(a)-5(c) illustrate this with an example of a job j that requires 3 datasets f_1, f_2 and f_3 that are replicated on data host sets $\{d_1, d_2\}$, $\{d_2, d_3\}$ and $\{d_1, d_4\}$ respectively. Figure 5(c) shows a possible dominating set for the graph of datasets and data hosts shown in Figure 5(b).

From Figure 5(b), we can build a reduced adjacency matrix A for graph G wherein $a_{ik} = 1$ if data host d_i contains f_k . Such an adjacency matrix is shown in Figure 6(a). Therefore, the problem of finding the minimum dominating sets for G is now equivalent to finding the sets of the least number of rows such that every column contains an entry of 1 in at least one of the rows. This problem has been studied extensively as the *Set Covering Problem* [19].

Christofides [20] provides an approximate tree search algorithm for finding a solution to the general Set Covering Problem. Based on this algorithm, we propose a mapping heuristic to find a minimum dominating set that ensures the smallest makespan. The heuristic is listed in Figure 7.

$$\begin{array}{c}
\begin{array}{ccc}
& f_1 & f_2 & f_3 \\
d_1 & \begin{pmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \\
d_2 \\
d_3 \\
d_4
\end{array} \\
\text{(a)}
\end{array}
\qquad
\begin{array}{c}
\begin{array}{ccc}
& f_1 & f_2 & f_3 \\
d_1 & \begin{pmatrix} 1 & 0 & 1 \\ 1 & 1 & 0 \\ - & - & - \\ d_2 & 1 & 1 & 0 \\ d_3 & 0 & 1 & 0 \\ - & - & - & - \\ d_1 & 1 & 0 & 1 \\ d_4 & 0 & 0 & 1 \end{pmatrix} \\
d_2 \\
d_3 \\
d_4
\end{array} \\
\text{(b)}
\end{array}
\end{array}$$

Fig. 6: (a) Adjacency Matrix. (b) Tableau.

At the start of the process, from the adjacency matrix, we create a tableau T consisting of K blocks of rows, where the k^{th} block consists of rows corresponding to data hosts that contain f_k . An example of a tableau generated from the adjacency matrix of Figure 6(a) is shown in Figure 6(b). The set of data hosts B keeps track of the current solution set of datahosts, the set E contains the datasets already covered by the solution set and the variable z keeps track of the makespan offered by the current solution set. The final solution set is stored in B_{final} . During execution, the blocks are searched sequentially starting from the k^{th} block where k is the smallest index, $1 \leq k \leq K$ such that $f_k \notin E$. Within the k^{th} block, let d_q^k mark the data host under consideration where q is a row pointer within block k . We add d_q^k to B and all the datasets for which the corresponding row contains 1, to E as they are already covered by d_q^k . These datasets are removed from consideration and the process then moves to the next uncovered block until $E = F_j$, that is, all the datasets have been covered. The function $MCT(B)$ computes the expected completion time for each compute resource combined with the solution set B and returns with the minimum of the completion time so found.

Through the recursive procedure outlined in the listing, the heuristic then backtracks and discovers other solution sets. The solution set that guarantees minimum makespan is then chosen as the final . The search terminates when the first block is exhausted. Therefore, before the tableau is created, we sort the rows of the adjacency matrix (that is, the data hosts) in the descending order of the number of columns with 1's (or the number of datasets contained). Also, in the tableau, the same sorting order is applied to the rows in each block. As the minimal dominating sets would obviously contain at least one of the datahosts with the maximum number of datasets, this increases the chances of more dominating sets being in the path of the search function within the proposed heuristic. Overall, the running time of the mapping heuristic is given by $O(MK^2)$ where MK^2 is the number of resource sets that are searched by the heuristic to find one that provides the least completion time.

Begin Main

1. For a job j , create the adjacency matrix A with data hosts forming the rows and datasets forming the columns.
2. Sort the rows of A in the descending order of the number of 1's in a row.
3. Create the tableau T from sorted A and begin with initial solution set $B_{final} = \phi$, $B = \phi$, $E = \phi$ and $z = \infty$
4. $\text{Search}(B_{final}, B, T, E, z)$
5. $S^j \leftarrow \{\{r\}, B_{final}\}$ where $r \in R$ such that $MCT(B_{final})$ is minimum

End Main

$\text{Search}(B_{final}, B, T, E, z)$

6. Find the minimum k , such that $f_k \notin E$. Let T_k be the block of rows in T corresponding to f_k . Set a pointer q to the top of T_k .
7. **while** q does not reach the end of T_k **do**
8. $F_T \leftarrow \{f_i | t_{qi} = 1, 1 \leq i \leq K\}$
9. $B \leftarrow B \cup \{d_q^k\}, E \leftarrow E \cup F_T$
10. **if** $E = F_j$ **then**
11. **if** $z > MCT(B)$ **then**
12. $B_{final} \leftarrow B, z \leftarrow MCT(B)$
13. **else** $\text{Search}(B_{final}, B, T, E, z)$
14. $B \leftarrow B - \{d_q^k\}, E \leftarrow E - F_T$
15. Increment q
16. **end**

$MCT(B)$

17. Find $r \in R$ such that the resource set $S^j = \{\{r\}, B\}$ gives minimum completion time

Fig. 7: Listing for SCP-based Mapping Heuristic.

Other heuristics that are possible or have been proposed include the ones described below:

Compute-First - In this mapping, a compute resource that ensures minimum computation time (T_c) is selected for the job first followed by choosing data hosts that have the best bandwidths to the selected resource. This is in contrast to our approach that places more importance on selection of data hosts. The running time of this heuristic is $O(MK)$.

Greedy - This heuristic builds the resource set by iterating through the list of datasets and making a greedy choice for the data host for accessing each dataset, followed by choosing the nearest compute resource for that data host. At the end of each iteration, it checks whether the compute resource so selected is better than the one selected in previous iteration when the data hosts selected previously are considered. This heuristic was presented in [14]. The running time of this heuristic is $O(KP)$.

Brute Force - In this case, all the possible resource sets for a particular job are generated and the one guaranteeing the MCT is chosen for the job. This is a generalisation of the heuristic presented in [11]. While this heuristic guarantees that the resource set selected will be the best for the job, it searches through MP^K resource sets at a time. This leads to unreasonably large search spaces for higher values of K . For example, for a job requiring 5 datasets with 20 possible data hosts and 20 available compute resources, the search space will consist of $(20 * 20^5) = 64 * 10^6$ resource sets.

A point to note is that the sets of datasets required by 2 or more jobs in the same set are not mutually exclusive. Any dataset that is transferred during from one resource to another is retained at the receiver and therefore, this presents an additional source of data to successive jobs requiring access to that dataset.

While the mapping heuristic finds a resource set for a single job, we wish to minimize the total *makespan* [10], the total time from the start of the scheduling to the completion of the last job, of the application consisting of N such data-intensive jobs. To that end, we apply the well-known MinMin heuristic, proposed in [10], to schedule the entire set of jobs. The MinMin heuristic submits the job with the smallest MCT to the compute resource (already selected by the mapping heuristic for the job) that guarantees it.

V. EXPERIMENTS

We have used GridSim with its new Data Grid capabilities [21] to simulate the data-intensive environment and evaluate the performance of scheduling algorithms. For evaluation, we have used the EU DataGrid topology based on the testbed shown in Figure 2. The details of the Grid resources used in our evaluation is shown in Table I. The resources in the actual testbed have gone through several configuration changes, not all of which are publicly available, and hence it was impossible to model their layout and CPU capability accurately. Instead, it was decided to create a configuration for each resource such that the modelled testbed in whole would reflect the heterogeneity of platforms and capabilities that is normally the characteristic of Grid resources. All the resources were simulated as clusters with a batch job management system using space-shared policy, as a front-end to single CPU processing nodes. The CPUs are rated in terms of MIPS (Million Instructions Per Sec). The resource at CERN was considered as a pure data source (data host) in our evaluation and hence, no jobs were submitted to it. To model resource contention caused by multiple users, we associate a *mean load* with each resource. The load factor is the ratio of the number of CPUs that are occupied to the total number of CPUs available within a resource. During the simulation, for each resource, we derive the instantaneous resource load from a Gaussian distribution with its mean as the load shown in Table I.

Similarly, we model the variability of the available network bandwidth by associating an availability factor with a link which is the ratio of the available bandwidth to the total bandwidth. During simulation, the instantaneous measure is derived from another Gaussian distribution centered around a mean availability factor assigned at random to each of the links.

Within this evaluation, we consider a universal set of datasets, each of which are replicated on one or more of the resources. Studies of similar environments [22] have shown that the size of the datasets follow a heavy-tailed distribution in which there are larger numbers of smaller size files and vice

TABLE I: Simulated configuration of resources within EDG testbed.

Resource Name (Location)	No. of Nodes	CPU Rating (MIPS)	Storage (TB)	Load
RAL (UK)	41	1140	2.75	0.9
Imperial College (UK)	52	1330	1.80	0.95
NorduGrid (Norway)	17	1176	1.00	0.9
NIKHEF (Netherlands)	18	1166	0.50	0.9
Lyon (France)	12	1320	1.35	0.8
CERN (Switzerland)	–	–	12	–
Milano (Italy)	7	1,000	0.35	0.5
Torino (Italy)	4	1330	0.10	0.5
Catania (Italy)	5	1200	0.25	0.6
Padova (Italy)	13	1,000	0.05	0.4
Bologna (Italy)	20	1140	5.00	0.8

versa. Therefore, we generate the set of datasets with sizes distributed according to the logarithmic distribution in the interval [1GB, 6GB]. The distribution of datasets depends on many factors itself including variations in popularity, the replication strategy employed and the nature of the fabric. Within our evaluation, we have used two commonly considered patterns of file distribution:

- *Uniform* : Here, the distribution of datasets is modeled on a uniform random probability distribution. In this scenario, each file is equally likely to be replicated at any site.
- *Zipf* : Zipf-like distributions follow a power law model in which the probability of occurrence of the i^{th} ranked file in a list of files is inversely proportional to i^{-a} where $a \leq 1$. In other words, a few files are distributed widely whereas most of files are found in one or two places. This models a scenario where the files are replicated on the basis of popularity. It has been shown that Zipf-like distributions holds true in cases such as requests for pages in World Wide Web where a few of the sites are visited the most [23]. This scenario has been evaluated for a Data Grid environment in related publications [24].

Henceforth, we will consider the distribution applied to be described by the variable *Dist*. We also control the distribution of datasets through a parameter called the *degree of replication* which is the maximum possible number of copies of any dataset in a Data Grid. The degree of replication in our evaluation is 5.

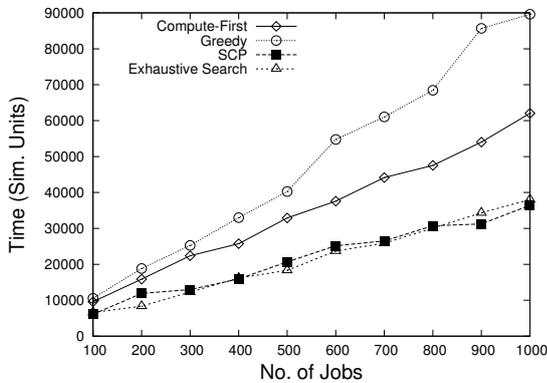
On the application side, there are three variables that determine the performance of the application: the size of the application or the number of jobs in the application (N), the number of datasets required by each job (K) and the computational size of a job ($Size(j)$) expressed in Million Instructions (MI). For each job, K datasets are selected at random from the universal set of datasets. For the purpose of comparison, we keep K a constant among all the jobs in a set although this is not a condition imposed on the heuristic itself. An experiment is described by the tuple $(N, K, Size, Dist)$

for which all the heuristics are evaluated. At the beginning of each experiment, the set of datasets, their distribution among the resources and the set of jobs are generated. This configuration is then kept constant while each of the four mapping heuristics are evaluated in turn. We have conducted 50 such experiments with different values for N , K , $Size$ and $Dist$ and in the next section, we present results of our evaluation.

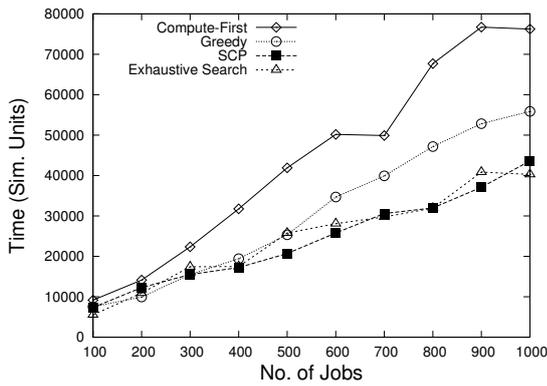
A. Results

TABLE II: Summary of Evaluations

Mapping Heuristic	Geometric Mean	Avg. deg.	Avg. rank
Compute-First	37593.71	69.01 (19.4)	3.63 (0.48)
Greedy	36927.44	71.86 (50.55)	3.23 (0.71)
SCP	24011.17	7.68 (10.42)	1.67 (0.6)
Brute Force	23218.49	3.87 (6.46)	1.47 (0.58)



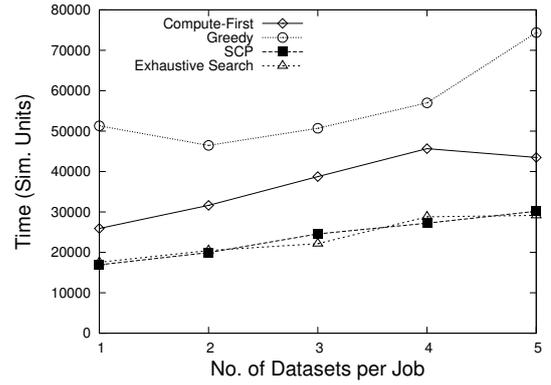
(a) Size=300000 MI,K=3, Dist=Uniform



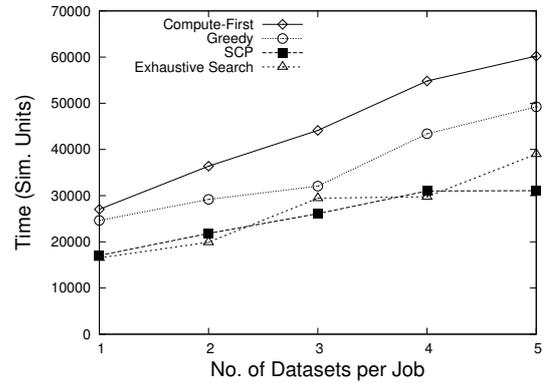
(b) Size=300000 MI,K=3, Dist=Zipf

Fig. 8: Makespan vs Number of Jobs.

The results of our evaluations are summarised in Table II and are based on the methodology provided in [9]. *SCP* refers to the heuristic proposed in this paper. For each mapping



(a) N=600,Size=300000 MI, Dist=Uniform

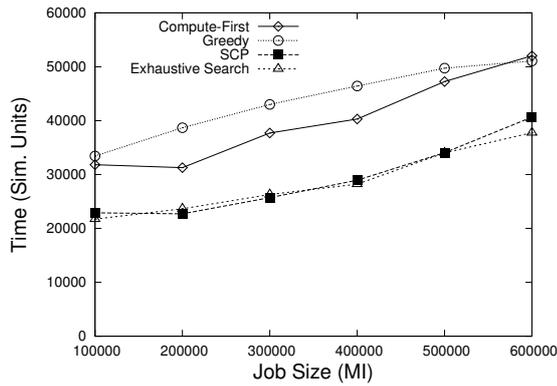


(b) N=600,Size=300000 MI, Dist=Zipf

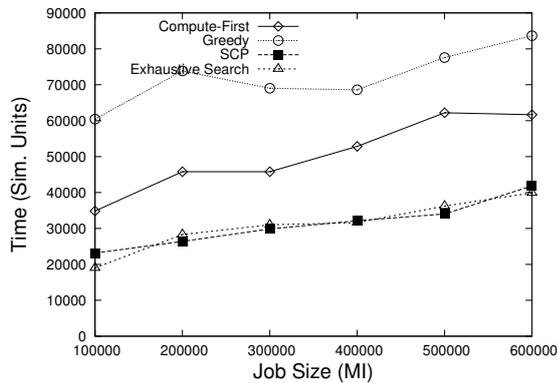
Fig. 9: Makespan vs Datasets per Job.

heuristic, the table contains three values: 1) *Geometric Mean* of the makespans, 2) *Average degradation (Avg. deg.)* from the best heuristic in an experiment and 3) *Average ranking (Avg. rank)* of each heuristic in an experiment. The geometric mean is used as the makespans vary in orders of magnitude according to parameters such as number of jobs per application set, number of files per job and the size of each job. Degradation for a heuristic is the difference between the makespan of that heuristic and that of the best heuristic for a particular experiment and expressed as a percentage of the latter. The average degradation is computed as an arithmetic mean over all experiments and the standard deviation of the population is given in the parantheses next to the means in the table. This is the measure of how far a heuristic is away from the best heuristic for an experiment. A lower number certainly means that the application is on an average the best one. The ranking is in the ascending order of makespans produced by the heuristics for each experiment, that is, lower the makespan, lower the rank of the heuristic. The standard deviation of the population is provided alongside the averages in the table.

The three values together provide a consolidated view of the



(a) N=600, K=3, Dist=Uniform



(b) N=600, K=5, Dist=Uniform

Fig. 10: Makespan vs Job Size.

performance of each heuristic. For example, we can see that on an average Compute-First and Greedy both perform worse than either SCP or Brute Force. However, the standard deviation of the population is much higher in the case of Greedy than that of Compute-First. Therefore, it can also be said that Compute-First can be expected to perform worst most of time. Indeed, in a few of the experiments, Greedy performed as good or even better than SCP while Compute-First never came close to the performance of the other heuristics.

As expected, between SCP and Brute Force, the latter is the clear winner having a consistently lower score than the former. However, the computational complexity of Brute Force means that as the number of datasets per job increases, the number of resource sets that need to be considered by the Brute Force heuristic increases dramatically. The geometric mean and average rank of SCP is close to that of Brute Force heuristic. The average rank is less than 2 for both heuristics which implies that in many scenarios, SCP provides a better performance than Brute Force.

This view is reinforced from the graphs in Figures 8-10 which show the effect of varying one of the variables, all

others kept constant. SCP and Brute-force give almost similar performance while either Greedy or Compute-First is the worst in almost all cases. The effect of job distribution is most visible on the Greedy heuristic. When the files are distributed according to the Zipf distribution, the performance of Greedy comes close to or in some cases, becomes as competitive as SCP. This is due to the fact that in Zipf distribution, there are most of the datasets are not replicated widely and therefore, there is not as much choice of data hosts as there is in Uniform distribution. In such a case, Greedy is able to form minimal resource sets. Also, it can be seen that as the number of jobs increases, the makespan of Compute-First and Greedy heuristic rise more steeply than the other two.

VI. CONCLUSION AND FUTURE WORK

We have presented the problem of mapping an application with a collection of jobs that require multiple datasets that are each replicated on multiple data hosts to Grid resources. We have also proposed a heuristic based on a solution to the Set Covering Problem. We have shown via simulation that the proposed heuristic is better than Compute-First and Greedy approaches and leads to schedules that are competitive with the exhaustive search (Brute Force) approach while being orders of magnitude faster.

As part of immediate future work, we plan to evaluate the SCP mapping heuristic using other task scheduling algorithms such as Max-min, Sufferage and Genetic Algorithms. The performance of the SCP-based heuristic in scenarios involving dependent tasks such as Directed Acyclic Graphs (DAGs) also needs to be investigated.

ACKNOWLEDGEMENT

We would like to thank Anthony Sulistio for his help with the use of GridSim. We would also like to thank Tianchi Ma, Kyong Hoon Kim and Chee Shin Yeo (GRIDS Lab), and the anonymous reviewers for their comments that have helped improve the quality of this paper.

REFERENCES

- [1] I. Foster and C. Kesselman, *The Grid: Blueprint for a Future Computing Infrastructure*. San Francisco, USA: Morgan Kaufmann Publishers, 1999.
- [2] A. Chervenak, I. Foster, C. Kesselman, C. Salisbury, and S. Tuecke, "The Data Grid: Towards an architecture for the distributed management and analysis of large scientific datasets," *Journal of Network and Computer Applications*, vol. 23, no. 3, pp. 187–200, 2000.
- [3] W. Hoschek, F. J. Jaen-Martinez, A. Samar, H. Stockinger, and K. Stockinger, "Data Management in an International Data Grid Project," in *Proceedings of the 1st IEEE/ACM International Workshop on Grid Computing (GRID '00)*. Bangalore, India: Springer-Verlag, Berlin, Germany, Dec. 2000.
- [4] N. Yamamoto, O. Tatebe, and S. Sekiguchi, "Parallel and Distributed Astronomical Data Analysis on Grid Datafarm," in *Proceedings of 5th IEEE/ACM International Workshop on Grid Computing (Grid 2004)*. Pittsburgh, USA: IEEE CS Press, Los Alamitos, CA, USA, Nov. 2004.
- [5] K. Ranganathan and I. Foster, "Decoupling Computation and Data Scheduling in Distributed Data-Intensive Applications," in *Proceedings of the 11th IEEE Symposium on High Performance Distributed Computing (HPDC)*. Edinburgh, UK: IEEE CS Press, Los Alamitos, CA, USA, July 2002.

- [6] W. H. Bell, D. G. Cameron, L. Capozza, A. P. Millar, K. Stockinger, and F. Zini, "Simulation of Dynamic Grid Replication Strategies in OporSim," in *Proceedings of the 3rd International Workshop on Grid Computing (GRID 02)*. Baltimore, MD, USA: Springer-Verlag, Berlin, Germany, 2002, pp. 46–57.
- [7] A. Takefusa, O. Tatebe, S. Matsuoka, and Y. Morita, "Performance Analysis of Scheduling and Replication Algorithms on Grid Datafarm Architecture for High-Energy Physics Applications," in *Proceedings of the 12th IEEE International Symposium on High Performance Distributed Computing (HPDC-12)*. Seattle, USA: IEEE CS Press, Los Alamitos, CA, USA, June 2003.
- [8] T. Phan, K. Ranganathan, and R. Sion, "Evolving toward the perfect schedule: Co-scheduling job assignments and data replication in wide-area systems using a genetic algorithm," in *Proceedings of the 11th Workshop on Job Scheduling Strategies for Parallel Processing*. Cambridge, MA: Springer-Verlag, Berlin, Germany, June 2005.
- [9] H. Casanova, A. Legrand, D. Zagorodnov, and F. Berman, "Heuristics for Scheduling Parameter Sweep Applications in Grid environments," in *Proceedings of the 9th Heterogeneous Computing Systems Workshop (HCW 2000)*. Cancun, Mexico: IEEE CS Press, Los Alamitos, CA, USA, 2000.
- [10] M. Maheswaran, S. Ali, H. J. Siegel, D. Hensgen, and R. F. Freund, "Dynamic Mapping of a Class of Independent Tasks onto Heterogeneous Computing Systems," *Journal of Parallel and Distributed Computing*, vol. 59, pp. 107–131, Nov 1999.
- [11] H. Mohamed and D. Epema, "An evaluation of the close-to-files processor and data co-allocation policy in multiclustes," in *Proceedings of the 2004 IEEE International Conference on Cluster Computing*. San Diego, CA, USA: IEEE CS Press, Los Alamitos, CA, USA, Sept. 2004.
- [12] A. Giersch, Y. Robert, and F. Vivien, "Scheduling tasks sharing files from distributed repositories," in *Proceedings of the 10th International Euro-Par Conference (EuroPar '04)*. Pisa, Italy: Springer-Verlag, Berlin, Germany, Sept. 2004.
- [13] G. Khanna, N. Vydyanathan, T. Kurc, U. Catalyurek, P. Wyckoff, J. Saltz, and P. Sadayappan, "A hypergraph partitioning-based approach for scheduling of tasks with batch-shared I/O," in *Proceedings of the 2005 IEEE International Symposium on Cluster Computing and the Grid (CCGrid 2005)*. Cardiff, UK: IEEE CS Press, May 2005.
- [14] S. Venugopal and R. Buyya, "A Deadline and Budget Constrained Scheduling Algorithm for e-Science Applications on Data Grids," in *Proceedings of the 6th International Conference on Algorithms and Architectures for Parallel Processing (ICA3PP-2005)*, ser. Lecture Notes in Computer Science, vol. 3719. Melbourne, Australia: Springer-Verlag, Berlin, Germany, Oct. 2005.
- [15] R. Gardner *et al.*, "The Grid2003 Production Grid: Principles and Practice," in *Proceedings of the 13th Symposium on High Performance Distributed Computing (HPDC 13)*. Honolulu, HI, USA: IEEE CS Press, Los Alamitos, CA, USA, June 2004.
- [16] A. Rajasekar, M. Wan, and R. Moore, "MySRB & SRB: Components of a Data Grid," in *Proceedings of the 11th IEEE International Symposium on High Performance Distributed Computing (HPDC-11)*. Edinburgh, UK: IEEE CS Press, Los Alamitos, CA, USA, 2002.
- [17] E. Seidel, G. Allen, A. Merzky, and J. Nabrzyski, "GridLab: a grid application toolkit and testbed," *Future Gener. Comput. Syst.*, vol. 18, no. 8, pp. 1143–1153, 2002.
- [18] S. Venugopal, R. Buyya, and L. Winton, "A Grid Service Broker for Scheduling Distributed Data-Oriented Applications on Global Grids," in *Proceedings of the 2nd Workshop on Middleware in Grid Computing (MGC 04)*. Toronto, Canada: ACM Press, New York, USA, Oct. 2004.
- [19] E. Balas and M. W. Padberg, "On the Set-Covering Problem," *Operations Research*, vol. 20, no. 6, pp. 1152–1161, 1972.
- [20] N. Christofides, *Graph Theory: An Algorithmic Approach*. Academic Publishers, London, UK, 1975, ch. Independent and Dominating Sets – The Set Covering Problem, pp. 30 – 57, ISBN 012 1743350 0.
- [21] A. Sulistio, U. Cibej, B. Robic, and R. Buyya, "A Tool for Modelling and Simulation of Data Grids with Integration of Data Storage, Replication and Analysis," University of Melbourne, Australia, Tech. Rep. GRIDS-TR-2005-13, Nov. 2005.
- [22] K. Park, G. Kim, and M. Crovella, "On the relationship between file sizes, transport protocols, and self-similar network traffic," in *Proceedings of the 1996 International Conference on Network Protocols (ICNP '96)*. Atlanta, GA, USA: IEEE CS Press, 1996.
- [23] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web caching and zipf-like distributions: evidence and implications," in *Proceedings of the 18th Annual Joint Conference of the IEEE Computer and Communications Societies (INFOCOM '99)*, New York, NY, USA, 1999.
- [24] D. G. Cameron, R. Carvajal-Schiaffino, A. P. Millar, C. Nicholson, K. Stockinger, and F. Zini, "Evaluating Scheduling and Replica Optimisation Strategies in OporSim," in *Proceedings of the 4th International Workshop on Grid Computing (Grid2003)*. Phoenix, AZ, USA: IEEE CS Press, Los Alamitos, CA, USA, Nov. 2003.