

# International Journal of High Performance Computing Applications

<http://hpc.sagepub.com>

---

## Steady-State Scheduling of Multiple Divisible Load Applications on Wide-Area Distributed Computing Platforms

L. Marchal, Y. Yang, H. Casanova and Y. Robert

*International Journal of High Performance Computing Applications* 2006; 20; 365

DOI: 10.1177/1094342006067474

The online version of this article can be found at:  
<http://hpc.sagepub.com/cgi/content/abstract/20/3/365>

---

Published by:

 SAGE Publications

<http://www.sagepublications.com>

Additional services and information for *International Journal of High Performance Computing Applications* can be found at:

**Email Alerts:** <http://hpc.sagepub.com/cgi/alerts>

**Subscriptions:** <http://hpc.sagepub.com/subscriptions>

**Reprints:** <http://www.sagepub.com/journalsReprints.nav>

**Permissions:** <http://www.sagepub.com/journalsPermissions.nav>

# STEADY-STATE SCHEDULING OF MULTIPLE DIVISIBLE LOAD APPLICATIONS ON WIDE-AREA DISTRIBUTED COMPUTING PLATFORMS

L. Marchal<sup>1</sup>  
Y. Yang<sup>2</sup>  
H. Casanova<sup>2,3</sup>  
Y. Robert<sup>1</sup>

## Abstract

Divisible load applications consist of an amount of data and associated computation that can be divided arbitrarily into any number of independent pieces. This model is a good approximation of many real-world scientific applications, lends itself to a natural master-worker implementation, and has thus received a lot of attention. The critical issue of divisible load scheduling has been studied extensively in previous work. However, only a few authors have explored the simultaneous scheduling of multiple such applications on a distributed computing platform. We focus on this increasingly relevant scenario and make the following contributions. We use a novel and more realistic platform model that captures some of the fundamental network properties of grid platforms. We formulate the steady-state multi-application scheduling problem as a linear program that expresses a notion of fairness between applications. This scheduling problem is NP-complete and we propose several heuristics that we evaluate and compare via extensive simulation experiments. Our main finding is that some of our heuristics can achieve performance close to the optimal and we quantify the trade-offs between achieved performance and heuristic complexity.

Key words: parallel computing, scheduling, divisible load, multiple applications, resource sharing

---

*The International Journal of High Performance Computing Applications*,  
Volume 20, No. 3, Fall 2006, pp. 365–381  
DOI: 10.1177/1094342006067474  
© 2006 SAGE Publications  
Figures 1–3 appear in color online: <http://hpc.sagepub.com>

## 1 Introduction

A *divisible load* application (Bharadwaj et al. 1996) consists of an amount of computation, or *load*, that can be divided into any number of independent pieces. This corresponds to a perfectly parallel job: any sub-task can itself be processed in parallel, and on any number of workers. The divisible load model is a good approximation for applications that consist of large numbers of identical, low-granularity computations, and has thus been applied to a wide spectrum of scientific problems in areas including image processing, volume rendering, bioinformatics, and even data mining. For further information on the model, we refer the reader to Bharadwaj, Ghose, and Robertazzi (2003), Robertazzi (2003, n.d.) and Ghose and Robertazzi (2003).

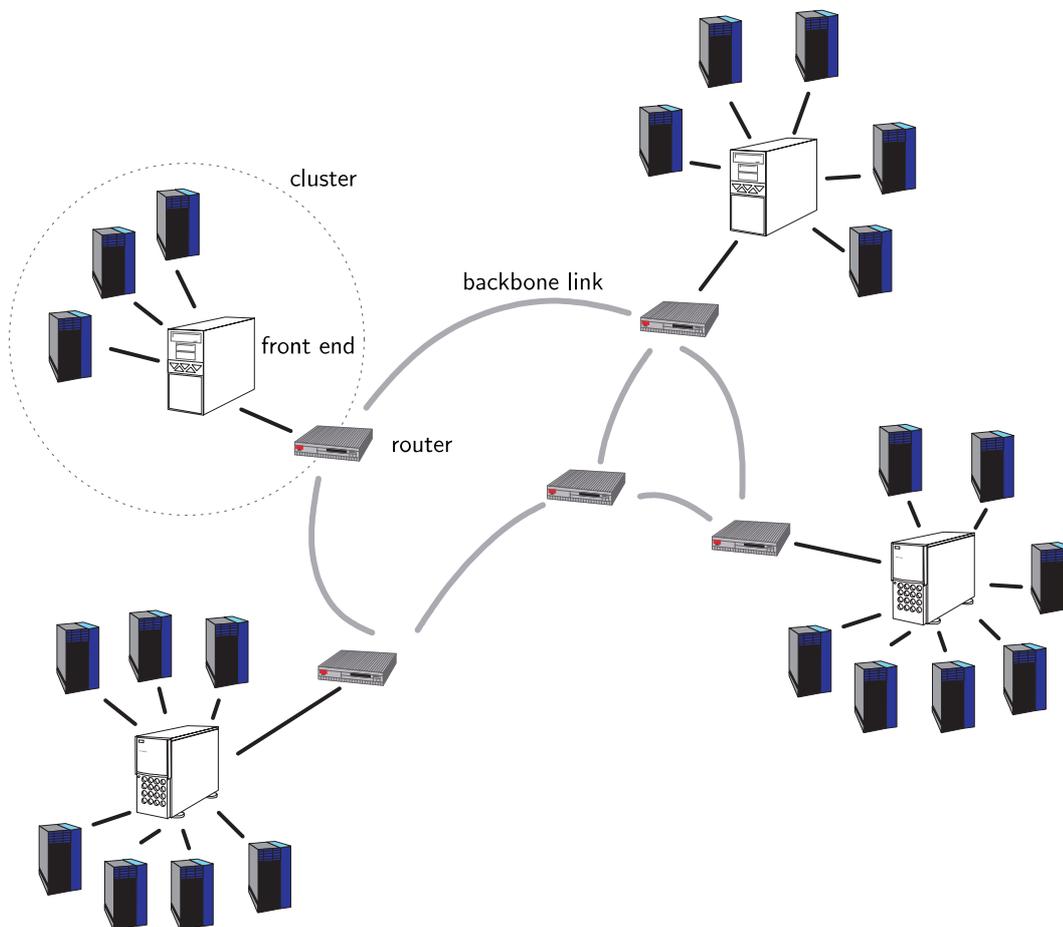
Divisible load applications are amenable to the simple master-worker programming model and can therefore be easily implemented and deployed on computing platforms ranging from small commodity clusters to computational grids. The main challenge, which has been studied extensively, is to *schedule* such applications effectively. However, large-scale platforms are not likely to be exploited in a mode dedicated to a single application. Furthermore, a significant portion of the mix of applications running on grid platforms are divisible load applications. At the extreme, a grid such as the CDF Analysis Farms (CAF) supports the concurrent executions of applications that are almost all divisible load applications. Therefore, it is critical to investigate the scheduling of *multiple* such applications that are executed simultaneously and compete for CPU and network resources.

A first analysis of the concurrent execution of multiple divisible load applications is provided in Bharadwaj and Barlas (2002). The authors target a simple platform composed of a bus network connecting a single master processor to a collection of heterogeneous worker processors. A more complex platform has been investigated in (Wong et al. 2003). In that paper, the authors introduce a (virtual) producer–consumer architecture where several masters (the sources of the multiple divisible loads) are fully connected to heterogeneous worker processors. The authors describe a strategy for balancing the total amount of work among the workers. Unfortunately, the results are mostly of theoretical interest as it is assumed that masters and workers can communicate with unlimited numbers of concur-

<sup>1</sup>LABORATOIRE LIP, CNRS-INRIA, ÉCOLE NORMALE SUPÉRIEURE DE LYON, FRANCE

<sup>2</sup>DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING, UNIVERSITY OF CALIFORNIA, SAN DIEGO, USA

<sup>3</sup>INFORMATION AND COMPUTER SCIENCES DEPARTMENT UNIVERSITY OF HAWAII AT MANOA, USA



**Fig. 1** Sample large-scale platform model.

rent messages, which is not likely to hold in practice. Yu and Robertazzi (2003), discuss how to apply divisible load theory to grid computing. They discuss job scheduling policies for a master-worker computation in which the workers are assumed to be only limited by their own network bandwidth and never by internet bandwidth. While possible, this assumption does not hold in general.

Our contributions are as follows: (i) we propose a new model for deploying and scheduling multiple divisible load applications on large-scale computing platforms, which is significantly more realistic than models used in previous work; (ii) we formulate a relevant multi-application steady-state divisible load scheduling problem, which is NP-complete; (iii) we propose several polynomial heuristics that

we evaluate and compare via extensive simulations. In our model, the target platform consists of a collection of clusters in geographically distributed institutions, interconnected via wide-area networks, as seen in Figure 1. The key benefit of this model is that it takes into account both the inherent hierarchy of the platform and the bandwidth-sharing properties of specific network links. In addition to the new platform model, we adopt a new scheduling objective. Rather than minimizing total application execution time (i.e. the “makespan”), our goal is to maximize the (priority-weighted) throughput in steady-state mode, i.e. the total load executed per time-period. There are three main reasons for focusing on the steady-state operation. First is *simplicity*, as steady-state scheduling is really a relaxation

of the makespan minimization problem in which the initialization and clean-up phases are ignored. One only needs to determine, for each participating resource, which fraction of time is spent computing for which application, and which fraction of time is spent communicating with which neighbor; the actual schedule then arises naturally from these quantities. Second is *efficiency*, as steady-state scheduling provides, by definition, a periodic schedule, which is described in compact form and can thus be implemented efficiently in practice. Third is *adaptability*: because the schedule is periodic, it is possible to dynamically record the observed performance during the current period, and to inject this information into the algorithm that will compute the optimal schedule for the next period. This makes it possible to react on the fly to resource availability variations.

In addition to the platform model and the scheduling objective described above, our approach enforces the constraint that the divisible load applications are processed fairly and allows for different application priorities.

## 2 Platform Model

Our platform model (see Figure 1) consists of a collection of clusters that are geographically distributed over the internet. Each cluster is equipped with a “front-end” processor (Bharadwaj et al. 1996), which is connected to a local router via a local-area link of limited capacity. These routers are used to connect each cluster to the internet. We model the interconnection of all the routers in our platform as a graph of internet backbone links.

The inter-cluster graph, denoted as  $\mathcal{G}_{ic} = (\mathcal{R}, \mathcal{B})$ , is composed of routers (the nodes in  $\mathcal{R}$ ) and backbone links (the edges in  $\mathcal{B}$ ). There are  $b = |\mathcal{B}|$  backbone links,  $l_1, \dots, l_b$ . For each link we have two parameters:  $bw(l_i)$ , the bandwidth available for a new connection, and  $max-connect(l_i)$ , the maximum number of connections (in both directions) that can be opened on this link by our applications. The model for the backbones is as follows. Each connection

is granted at most a fixed amount of bandwidth equal to  $bw(l_i)$ , up to the point where a maximum number of connections are simultaneously opened, at which point no more connections can be added. This model is justified by the bandwidth-sharing properties observed on wide-area links: when such a link is a bottleneck for an end-to-end TCP flow, several extra flows can generally be opened on the same path and they each receive the same amount of bandwidth as the original flow. This behavior can be due to TCP itself (e.g. congestion windows), or to the fact that the number of flows belonging to a single application is typically insignificant when compared with the total number of flows going through these links. This property is often exploited by explicitly opening parallel TCP connections (e.g. in the GridFTP project; Allcock et al. 2001) and we have observed it in our own measurements (Casanova 2004). The constraint imposed on the number of allowed connections makes it possible to limit the network usage of applications, which is a likely requirement for future production grid platforms with many applications and users competing for resources.

Compute resources consist of  $K$  clusters  $C^k$ ,  $1 \leq k \leq K$ . In full generality, we should represent each  $C^k$  as a node-weighted, edge-weighted graph  $G_k = (V_k, E_k)$ , but we simplify the model. For each cluster  $C^k$ , we only retain  $C_{master}^k$ , the front-end processor, which is connected to  $C_{router}^k$ , one of the routers in  $\mathcal{R}$ . The idea is that  $C_{master}^k$  represents the cumulated power of the computing resources in cluster  $C^k$  (as shown in Figure 2). This amounts to assuming that the architecture of the cluster is a star-shaped network, whose center is the front-end processor  $C_{master}^k$ . It is known that, for the purpose of running divisible load applications,  $C_{master}^k$  and the leaf processors are together “equivalent” to a single processor whose speed  $s_k$  can be determined as in Robertazzi (1993), Bataineh, Hsiung, and Robertazzi (1994), and Banino et al. (2004). In fact, it has also been shown that a tree topology is equivalent to a single processor (Bataineh, Hsiung, and Robertazzi

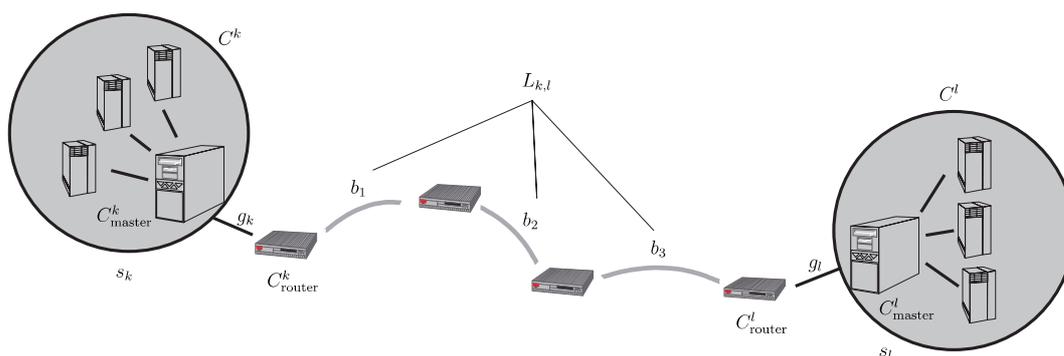


Fig. 2 Notations for the platform model.

1994; Barlas 1998; Beaumont et al. 2003), and thus our model encompasses cases in which the local-area network in each institution is structured as a tree. Consequently, we need only two parameters to characterize each cluster:  $s_k$ , the cumulated speed of  $C^k$  including  $C_{\text{master}}^k$  and the cluster's processors, and  $g_k$ , the bandwidth of the link connecting  $C_{\text{master}}^k$  to  $C_{\text{master}}^k$ . This link is modeled as follows: any number of connections may share the link, but they each receive a portion of the total available bandwidth, and the sum of these portions cannot exceed  $g_k$ , which is known to be a reasonable model for local-area links. Note that this link may correspond to several local area physical links.

Finally, we assume that the routing between clusters is fixed. The routing table contains an ordered list  $L_{k,l}$  of backbone links for a connection from cluster  $C^k$  to cluster  $C^l$ , i.e. from router  $C_{\text{router}}^k$  to router  $C_{\text{router}}^l$ . As shown in Figure 2, some intermediate routers may not be associated to any cluster. Also, no specific assumption is made on the interconnection graph. Our model uses realistic bandwidth assignments: we determine the bottleneck link for each end-to-end connection and use the bandwidth-sharing properties of this link (either local-area or backbone) to determine the amount of bandwidth allocated to each connection.

To the best of our knowledge, this model is the first attempt at modeling relatively complex network topologies along with realistic bandwidth-sharing properties for the purpose of large-scale application scheduling research. We contend that this model provides a major first step in the development of application-level scheduling strategies that are truly relevant to the new class of platforms brought about by grid infrastructures.

### 3 Steady-State Scheduling of Multiple Applications

The steady-state approach was pioneered by Bertsimas and Gamarnik (1999). Steady-state scheduling allows the relaxation of the scheduling problem in many ways. Indeed, initialization and clean-up phases are ignored, and the emphasis is on the design of a *periodic* schedule. Precise ordering and allocation of tasks and messages are not required. The key idea is to characterize the activity of each resource during each time-unit: which (rational) fraction of time is spent computing for which application, which (rational) fraction of time is spent receiving or sending to which neighbor. Such activity variables are used to construct a linear program that characterizes the global behavior of the system. Once each activity variable has been computed, the periodic schedule is known: we simply scale the rational values to obtain integer numbers, and the length of the period of the schedule is determined by this scaling. We outline below the construction step-by-step.

#### 3.1 Steady-State Equations

We consider  $K$  divisible load applications,  $A_k$ ,  $1 \leq k \leq K$ , with cluster  $C^k$  initially holding all the input data necessary for application  $A_k$ . For each application we define a "priority factor",  $\pi_k$ , that quantifies its relative worth. For instance, computing two units of load per time unit for an application with priority factor 2 is as worthwhile/profitable as computing one unit of load for an application with priority factor 1. This concept makes it possible to implement notions of application priorities for resource sharing. We can easily refine the priority model and define  $\pi_{k,l}$  as the priority factor to execute a fraction of application  $A_k$  onto cluster  $C^l$ . Similarly, our method is easily extensible to the case in which more than one application originates from the same cluster. We start with the following three definitions:

$w_k$  and  $\delta_k$  (**load unit size for  $A_k$** ) – The divisible applications may be of different types. For instance one application may deal with files and another with matrices. We divide each application into load units (a file, or a matrix). We let  $w_k$  be the amount of computation required to process a load unit for application  $A_k$ . Similarly,  $\delta_k$  is the size (in bytes) of a load unit for application  $A_k$ .

$\alpha_{k,l}$  (**fraction of  $A_k$  executed by  $C^l$** ) – Each cluster  $C^k$  initially holds input data for application  $A_k$ . Within a time-unit,  $C^k$  will devote a fraction of the time to process load units for application  $A_k$ . But cluster  $C^k$  can also be used to process loads that originate from another cluster  $C^l$ , i.e. from application  $A_l$ . Reciprocally, portions of application  $A_k$  may be executed by other clusters. We let  $\alpha_{k,l}$  be the portion of load for application  $A_k$  that is sent by  $C^k$  and computed on cluster  $C^l$  within a time-unit.  $\alpha_{k,k}$  denotes the portion of application  $A_k$  that is executed on the local cluster.

$\beta_{k,l}$  (**connections from  $C^k$  to  $C^l$** ) – Cluster  $C^k$  opens  $\beta_{k,l}$  network connections to send the portion  $\alpha_{k,l}$  of application  $A_k$  that is destined to cluster  $C^l$ .

We can reformulate the application model in terms of multiple jobs. Each job is composed of (an infinite number of) independent tasks (e.g. files to process or matrices to invert). The parameters  $w_k$  and  $\delta_k$  represent the task granularity of application  $A_k$ . In reality, only integer numbers of tasks should be considered. Here we relax the problem and deal with *rational* number of tasks in the steady-state equations that follow. Allowing rational values amounts to freely varying the granularity, thereby processing each application as if it was truly divisible. Later in Section 3.2 we briefly explain how to retrieve an effective schedule with an integer number of tasks for each application.

With the above definitions, it takes  $\frac{\alpha_{k,l} \cdot w_k}{s_l}$  time-units to process  $\alpha_{k,l}$  load units of application  $A_k$  on cluster  $C^l$ . Similarly, it takes  $\frac{\alpha_{k,l} \cdot \delta_k}{g_{k,l}}$  time-units to send  $\alpha_{k,l}$  load units of application  $A_k$  along a single network connection from router  $C_{router}^k$  to router  $C_{router}^l$ , where  $g_{k,l}$  is the minimum bandwidth available for one connection on a route from cluster  $C^k$  to cluster  $C^l$ , i.e.  $g_{k,l} = \min_{l_i \in L_{k,l}} \{bw(l_i)\}$ .

The first steady-state equation states that a cluster  $C^k$  cannot compute more load units per time unit than allowed by its speed  $s_k$ :

$$\forall C^k, \quad \sum_l \alpha_{l,k} \cdot w_l \leq s_k \quad (1)$$

With steady-state scheduling we do not need to determine the precise ordering in which the different load types are executed by  $C^k$ : instead we take a macroscopic point of view and simply bound the total amount of load processed every time-unit.

The second steady-state equation bounds the amount of load that requires the use of the serial link between cluster  $C^k$  and the external world, i.e. between  $C_{master}^k$  and  $C_{router}^k$ :

$$\forall C^k, \quad \underbrace{\sum_{l \neq k} \alpha_{k,l} \cdot \delta_k}_{\text{(outgoing data)}} + \underbrace{\sum_{j \neq k} \alpha_{j,k} \cdot \delta_j}_{\text{(incoming data)}} \leq g_k \quad (2)$$

This equation states that the available bandwidth  $g_k$  is not exceeded by the requirements of all the traffic outgoing from and incoming to cluster  $C^k$ . Again, there is no need to specify the precise ordering of the communications along the link. Note that we assume that the time to execute a portion of an application's load, or to communicate it along a serial link, is proportional to its size in number of load units: this amounts to fixing the granularity and to manipulating load units. Start-up costs could be included in the formulas, but at the price of technical difficulties: only asymptotic performance can be assessed in that case (Beaumont, Legrand, and Robert 2003).

Next we must bound the utilization of the backbone links. Our third equation states that on each backbone link  $l_i$ , there should be no more than  $max-connect(l_i)$  opened connections:

$$\sum_{\{k,l\}, l_i \in L_{k,l}} \beta_{k,l} \leq max-connect(l_i) \quad (3)$$

The fourth equation states that there is enough bandwidth on each path from a cluster  $C^k$  to a cluster  $C^l$ :

$$\alpha_{k,l} \cdot \delta_k \leq \beta_{k,l} \times g_{k,l}. \quad (4)$$

The last term  $g_{k,l}$  in Equation 4 was defined earlier as the bandwidth allotted to a connection from  $C^k$  to  $C^l$ . This bandwidth is simply the minimum of the  $bw(l_i)$ , taken over all links  $l_i$  that constitute the routing path from  $C^k$  to  $C^l$ . We multiply this bandwidth by the number of opened connections to derive the constraint on  $\alpha_{k,l}$ .

Finally there remains the definition of an optimization criterion. Let  $\alpha_k = \sum_{l=1}^K \alpha_{k,l}$  be the load processed for application  $A_k$  per time unit. To achieve a fair balance of resource allocations one could execute the same number of load units per application, and try to maximize this number. However, some applications may have higher priorities than others, hence the introduction of the priority factors  $\pi_k$  in the objective function:

$$\text{MAXIMIZE } \min_k \left\{ \frac{\alpha_k}{\pi_k} \right\}. \quad (5)$$

This maximization corresponds to the well-known MAX-MIN fairness strategy (Bertsekas and Gallager 1987) between the different loads, with coefficients  $1/\pi_k$ ,  $1 \leq k \leq K$ . The constraints and the objective function form a linear program:

$$\text{MAXIMIZE } \min_k \left\{ \frac{\alpha_k}{\pi_k} \right\}$$

UNDER THE CONSTRAINTS

$$\left\{ \begin{array}{ll} (6a) & \forall C^k, \quad \sum_l \alpha_{k,l} = \alpha_k \\ (6b) & \forall C^k, \quad \sum_l \alpha_{l,k} \cdot w_l \leq s_k \\ (6c) & \forall C^k, \quad \sum_{l \neq k} \alpha_{k,l} \cdot \delta_k + \sum_{j \neq k} \alpha_{j,k} \cdot \delta_j \leq g_k \\ (6d) & \forall i, \quad \sum_{L_{k,l} \ni l_i} \beta_{k,l} \leq max-connect(l_i) \\ (6e) & \forall k, l, \quad \alpha_{k,l} \cdot \delta_k \leq \beta_{k,l} \cdot g_{k,l} \\ (6f) & \forall k, l, \quad \alpha_{k,l} \geq 0 \\ (6g) & \forall k, l, \quad \beta_{k,l} \in \mathbb{N} \end{array} \right. \quad (6)$$

This program is mixed as the  $\alpha_{k,l}$  are rational numbers but the  $\beta_{k,l}$  are integers. Given a platform  $\mathcal{P}$  and computational priorities  $(\pi_1, \dots, \pi_K)$ , we define a *valid allocation* for the steady-state mode as a set of values  $(\alpha, \beta)$  such that Equations (6) are satisfied. Since this program involves integer variables there is little hope that an optimal solution could be computed in polynomial time. It turns out that this is an NP-hard problem, as shown in Section 4. However, the program captures all the con-

straints to be satisfied, and we can reconstruct a periodic schedule for every valid allocation (see Section 3.2). We derive several heuristics to compute valid allocations in Section 5, and we assess their performances in Section 6.

### 3.2 Reconstructing a Periodic Schedule

Once one has obtained a solution to the linear program defined in the previous section, say  $(\alpha, \beta)$ , one needs to reconstruct a (periodic) schedule, that is a way to decide in which specific activities each computation and communication resource is involved during each period. This is straightforward because the divisible load applications are independent of each other. We express all the rational numbers  $\alpha_{k,l}$  as  $\alpha_{k,l} = \frac{u_{k,l}}{v_{k,l}}$ , where the  $u_{k,l}$  and the  $v_{k,l}$  are relatively prime integers. The period of the schedule is set to  $T_p = \text{lcm}_{k,l}(v_{k,l})$ . In steady-state, during each period of duration  $T_p$ :

- Cluster  $C^k$  computes, for each non-zero value of  $\alpha_{l,k}$ ,  $\alpha_{l,k} \cdot T_p$  load units of application  $A_l$ . If  $l = k$  the data is local, and if  $k \neq l$ , the data corresponding to this load has been received during the previous period. These computations are executed in any order. Equation 2 ensures that  $\frac{\sum_l \alpha_{l,k} \cdot w_l \cdot T_p}{T_p} \leq s_k$ , hence  $C^k$  can process all its load.
- Cluster  $C^k$  sends, for each non-zero value of  $\alpha_{k,l}$ ,  $\alpha_{k,l} \cdot \delta_k \cdot T_p$  load units of application  $A_k$ , to be processed by cluster  $C^l$  during the next period. Similarly, it receives, for each non-zero value of  $\alpha_{j,k}$ ,  $\alpha_{j,k} \cdot z_{c_j} \cdot T_p$  load units for application  $A_j$ , to be processed locally during the next period. All these communications share the serial link, but Equation 2 ensures that  $\frac{\sum_{l \neq k} \alpha_{k,l} \cdot \delta_k \cdot T_p + \sum_{j \neq k} \alpha_{j,k} \cdot \delta_j \cdot T_p}{T_p} \leq g_k$ , hence the link bandwidth is not exceeded.

Obviously, the first and last period are different: no computation takes place during the first period, and no communication during the last one. Altogether, we have a periodic schedule, which is described in compact form: we have a polynomial number of intervals during which each processor is assigned a given load for a prescribed application.

### 4 Complexity

In this section we establish a complexity result: optimizing the throughput is NP-hard. We start with the formulation of the associated decision problem, and we proceed to the proof. Note that we cannot use a straightforward

reduction from a multicommodity flow problem such as problem ND47 in (Ausiello et al. 1999), because there is no prescribed location where each work should be executed.

**Definition 1. (STEADY-STATE-DIVISIBLE-LOAD( $\mathcal{P}$ ,  $\pi$ ,  $\rho$ )).** Given a platform  $\mathcal{P}$  and a set  $\mathcal{A}$  of divisible application with priority factors  $(\pi_1, \dots, \pi_K)$  and a throughput bound  $\rho$ , is there a valid allocation  $(\alpha, \beta)$  such that  $\min_k \left\{ \frac{\alpha_k}{\pi_k} \right\} \geq \rho$  ?

**Theorem 1. STEADY-STATE-DIVISIBLE-LOAD( $\mathcal{P}$ ,  $\pi$ ,  $\rho$ ) is NP-complete.**

**Proof.** The proof is detailed in the Appendix, at the end of the paper.

## 5 Heuristics

We propose several heuristics to solve our scheduling problem. We first propose a greedy heuristic, and then heuristics that are based on the rational solution to the mixed linear program derived in Section 3.

### 5.1 Greedy Heuristic

Our greedy heuristic, which we simply call G, allocates resources to one of the  $K$  applications in a sequence of steps. More specifically, at each step the heuristic (i) selects an application  $A_k$ ; (ii) determines on which cluster  $C^l$  the work will be executed (locally if  $l = k$ , on some remote cluster otherwise); and (iii) decides how much work to execute for this application. The intuition for how these choices can be made is as follows:

- One should select the application that has received the smallest relative share of the resource so far, that is the one for which  $\alpha_k/\pi_k$  is minimum, where  $\alpha_k = \sum_l \alpha_{k,l}$ . Initially,  $\alpha_k = 0$  for all  $k$ , so one can break ties by giving priority to the application with the highest priority factor  $\pi_k$ .
- Compare the pay-off of computing on the local cluster with the pay-off of opening one connection to each remote cluster. Choose the most profitable cluster, say  $C^l$ .
- Allocate an amount of work that does not overload  $C^l$  so that it will not be usable by other applications.

Again, let  $g_{k,l} = \min_{l_i \in L_{k,l}} \{bw(l_i)\}$  be the minimum bandwidth available for one connection on a route from cluster  $C^k$  to cluster  $C^l$ . The greedy heuristic, which we denote by G, is formalized as follows:

1. Let  $L = \{C^1, \dots, C^K\}$ . Initialize all  $\alpha_{k,l}$  and  $\beta_{k,l}$  to 0.

2. If  $L$  is empty, exit.
3. **Select application** – Sort  $L$  by non-decreasing values of  $\left(\frac{\alpha_k}{\pi_k}\right)$ . Break ties by choosing the application with higher priority first. Let  $k$  be the index of the first element of  $L$ . Select  $A_k$ .
4. **Select cluster** – For each cluster  $C^m$  where  $m \neq k$ , compute the work (i.e. number of load units for  $A_k$ ) that can be executed using a single connection:  $\text{benefit}_m = \min\left\{\frac{g_k}{\delta_k}, \frac{g_{k,m}}{\delta_k}, \frac{g_m}{\delta_k}, \frac{s_m}{w_k}\right\}$ . Locally, one can achieve  $\text{benefit}_k = \left(\frac{s_k}{w_k}\right)$ . Select  $C^l$ ,  $1 \leq l \leq K$  so that  $\text{benefit}_l$  is maximal. If  $\text{benefit}_l = 0$  (i.e. no more work can be executed), then remove  $C^k$  from list  $L$  and go to step 2.
5. **Determine amount of work** – If  $k \neq l$  (remote computation), allocate  $\text{alloc} = \text{benefit}_l$  units of load to cluster  $C^l$ . If  $k = l$  (local computation), allocate only  $\text{alloc} = \max_{m \neq k} \left\{ \min\left\{ \frac{g_k}{\delta_k}, \frac{g_{k,m}}{\delta_k}, \frac{g_m}{\delta_k}, \frac{s_m}{w_k} \right\} \right\}$  units of load. This last quantity is the largest amount that could have been executed on  $C^k$  for another application and is used to prevent over-utilization of the local cluster early on in the scheduling process.
6. **Update variables** –
  - Decrement speed of target cluster  $C^l$ :  $s_l \leftarrow s_l - \text{alloc} \cdot w_k$
  - Allocate work:  $\alpha_{k,l} \leftarrow \alpha_{k,l} + \text{alloc}$
  - In case of a remote computation (if  $k \neq l$ ) update network characteristics:

$$\begin{aligned} \forall l_i \in L_{k,l}, \\ \text{max-connect}(l_i) &\leftarrow \text{max-connect}(l_i) - 1 \\ g_k &\leftarrow g_k - \text{alloc} \cdot \delta_k, \quad g_l \leftarrow g_l - \text{alloc} \cdot \delta_k, \\ \beta_{k,l} &\leftarrow \beta_{k,l} + 1 \end{aligned}$$

7. Go to step 2.

## 5.2 LP-Based Heuristics

The linear program given in Section 3 is a mixed integer/rational numbers linear program since the variables  $\beta_{k,l}$  take integer values and variables  $\alpha_{k,l}$  may be rational. This mixed LP (MLP) formulation gives an exact optimal solution to the scheduling problem, while a rational LP formulation allows rational  $\beta_{k,l}$  and gives an *upper bound* of the optimal solution. As solving a mixed linear program is known to be hard, we propose several heuristics

based on the relaxation of the problem: we first solve the linear program over the rational numbers with a standard method (e.g. the Simplex algorithm). We then try to derive a solution with integer  $\beta_{k,l}$  from the rational solution.

**5.2.1 LPR: round-off** The most straightforward approach is to simply round rational  $\beta_{k,l}$  values to the largest smaller integer. Formally, if  $(\tilde{\alpha}_{k,l}, \tilde{\beta}_{k,l})$  is a rational solution to the linear program, we build the following solution:

$$\forall k, l, \quad \hat{\beta}_{k,l} = \lfloor \tilde{\beta}_{k,l} \rfloor, \quad \hat{\alpha}_{k,l} = \min\left\{ \tilde{\alpha}_{k,l}, \frac{\lfloor \tilde{\beta}_{k,l} \rfloor \cdot g_{k,l}}{\delta_k} \right\}.$$

With these new values, we have  $\hat{\beta}_{k,l} \leq \tilde{\beta}_{k,l}$  and  $\hat{\alpha}_{k,l} \leq \tilde{\alpha}_{k,l}$  for all indices  $k, l$ . Furthermore,  $(\hat{\alpha}, \hat{\beta})$  is a valid solution to the mixed linear program (6) in which all  $\hat{\beta}_{k,l}$  take integer values. We label this method LPR.

**5.2.2 LPRG: round-off + greedy** Rounding down all the  $\beta_{k,l}$  variables with LPR may lead to a very poor result as the remaining network capacity is unutilized. The LPRG heuristic reclaims this residual capacity by applying the technique described in Section 5.1. Intuitively, LPR gives the basic framework of the solution, while the greedy heuristic refines it.

**5.2.3 LPRR: randomized round-off** Relaxing an integer linear program into rational numbers is a classical approach, and several solutions have been proposed. Among them is the use of randomized approximation. Motwani, Naor, and Raghavan (1996) propose this approach to solve a related problem, the multicommodity flow problem. Using Chernoff bounds, they prove that their algorithm leads, with high probability, to a feasible solution that achieves the optimal throughput. Although this theoretical result seems attractive, it has some drawbacks for our purpose. First, our problem is not a multicommodity flow problem: instead of specifying a set of flow capacities for between node pairs, we have global demands for the sum of all flows leaving each node (representing the total amount of work sent by this node). Second, to obtain their optimality result, the authors in Motwani et al. (1996) rely on the assumption that the capacity of each edge is not smaller than a bound ( $5.2 \times \ln(4m)$  where  $m$  is the number of edges), and we do not have a similar property here. Third, there are two cases of failure in the randomized algorithm (even though the probability of such failures is proved to be small): either the algorithm provides a solution whose objective func-

tion is suboptimal (which is acceptable), or it provides a solution which does not satisfy all the constraints (which is not acceptable).

Coudert and Rivano (2002) proposed a rounding heuristic based on the method of Hochbaum (1997, Chapter 11) in the context of optical networks. Their method seems more practical as it always provides a feasible solution. We use a similar approach and our heuristic, LPRR, works as follows:

1. Solve the original linear program with rational numbers. Let  $(\tilde{\alpha}_{k,l}, \tilde{\beta}_{k,l})$  be the solution.
2. Choose a route  $k, l$  at random, such that  $\tilde{\beta}_{k,l} \neq 0$ .
3. Randomly choose  $X_{k,l} \in \{0, 1\}$  with probability  $P(X_{k,l} = 1) = \tilde{\beta}_{k,l} - \lfloor \tilde{\beta}_{k,l} \rfloor$ .
4. Assign the value  $v = \lfloor \tilde{\beta}_{k,l} \rfloor + X$  to  $\beta_{k,l}$  by adding the constraint  $\beta_{k,l} = v$  to the linear program.
5. If there is at least a route  $k, l$  for which no  $\beta_{k,l}$  value has been assigned, go to step 2.

Note that LPRR solves  $K^2$  linear programs, and is thus much more computationally expensive than our other LP-based heuristics.

## 6 Experimental Results

### 6.1 Methodology

In this section, we use simulation to evaluate the G, LPR, LPRG, and LPRR heuristics described in Section 5. Ideally the objective values achieved by these heuristics should be compared with the optimal solution, i.e. the solution to the mixed linear problem. However, solving the mixed linear problem takes exponential time and we cannot compute its solution in practice. Instead we use the solution to the *rational* linear problem as a comparator, as it provides an upper bound on the optimal solution (i.e. it may not be achievable in practice as  $\beta_{k,l}$  values must be integers).

One important question for creating relevant instances of our problem is that of the network topology. Indeed, the underlying topology impacts the relative performance of different heuristics. We evaluated our heuristics on two classes of network topologies. First, we generated a comprehensive set of topologies with Tiers (Calvert, Doar, and Zegura 1997), which has been widely accepted as a generator of realistic wide-area network topologies. Second, to try to get more insight into whether the nature of the network topology has a large impact on the relative performance of our heuristics, we generated simple random graphs in which each pair of nodes was connected

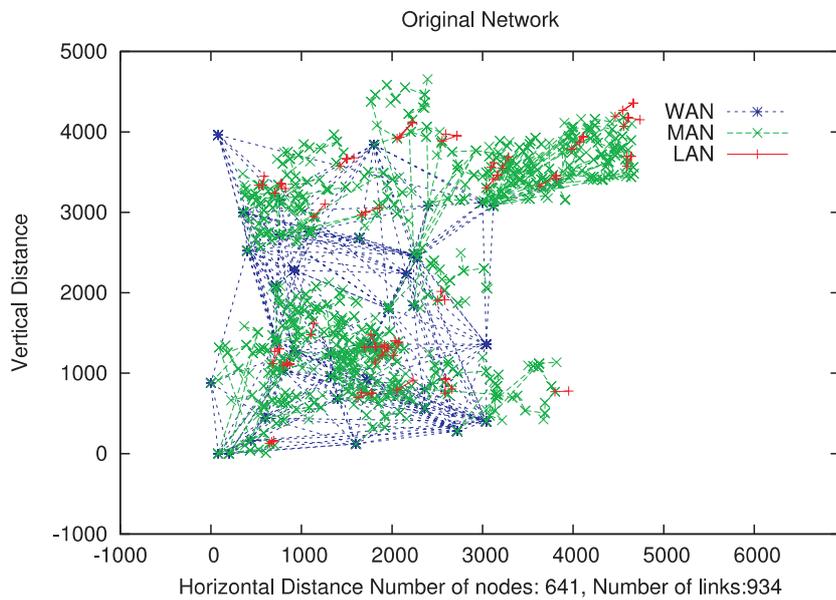
with a certain probability. We describe both topologies in detail below.

**(1) Tiers-generated topologies.** We randomly generated 100 two-level topologies with Tiers, each topology containing 40 WAN nodes and 30 MAN networks each containing 20 MAN nodes. We did not generate LAN networks as in our model we abstract them as a single cluster/site that delivers computation to the applications. We set a high connection redundancy value to reflect the rich connectivity between backbone nodes. Each of these topologies contains approximately 700 nodes. For each topology, we randomly select  $K = 5, 7, \dots, 90$  nodes as clusters participating in the computation of divisible load applications. For these  $K$  nodes we determine all pair-wise shortest paths (in hops), and we then delete the nodes not on any shortest paths, so as to be left with the topology interconnecting the sites participating in the computation. Note that in this “pruned” topology, there are nodes that we did not originally select but happen to be on the shortest paths between nodes that we had selected. We consider these nodes purely as routers that do not perform any computation, which can be easily expressed in the linear program defined in Section 3.1 by adding corresponding constraints but not modifying the objective function. Figure 3(a) shows a sample original Tiers topology, and Figure 3(b) shows the corresponding pruned topology.

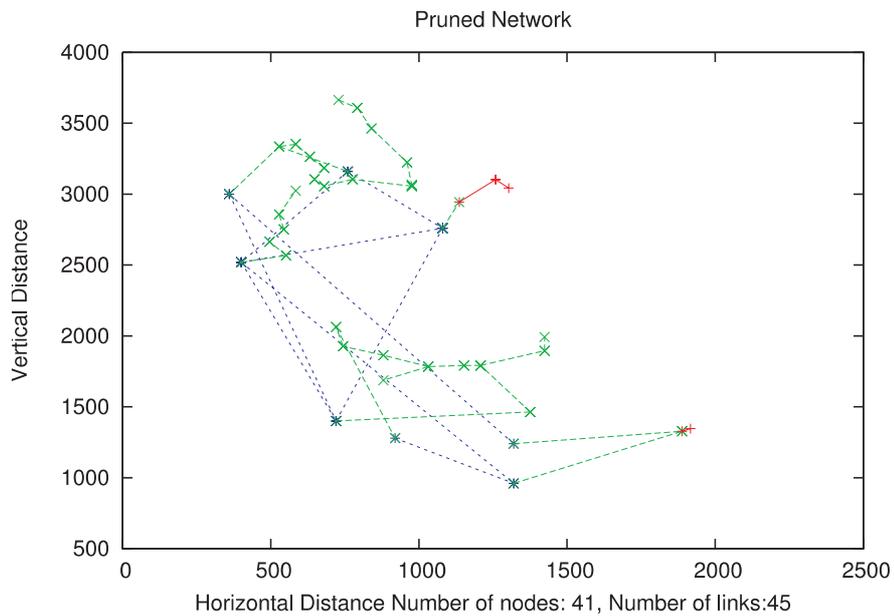
Once the topology is specified, we assign ranges of values to  $s_k, g_k, \text{max-connect}(l_k), \delta_k, w_k,$  and  $\pi_k$  as follows. The local bandwidth at each site,  $g_k$ , and the link bandwidth,  $bw(l_i)$ , is set according to a comprehensive measurement of internet end-to-end bandwidths (Lee and Stepanek 2001). This study shows that the logarithm of observed data transfer rates are approximately normally distributed with mean  $\log(2000 \text{ kbits/sec})$ , and standard deviation  $\log(10)$ , which we use to generate random values for  $g_k$  and  $bw(l_i)$ . We generate all the other parameters according to uniform distributions with ranges shown in Table 1. For each of our pruned Tiers topologies we generate 10 platform configurations with random instantiations of the

**Table 1**  
**Platform parameters for Tiers-generated topologies**

Parameter	Distribution
$K$	5, 7, ..., 90
$\log(bw(l_k)), \log(g_k)$	normal ( $mean = \log(2000)$ , $std = \log(10)$ )
$s_k$	uniform, 1000–10000
$\text{max-connect}, \delta_k, w_k, \pi_k$	uniform, 1–10



(a) Sample Tiers topology – 700 nodes.



(b) Sample pruned Tiers topology – 45 nodes.

**Fig. 3** Sample full and pruned Tiers topology.

**Table 2**  
Platform parameters for random graph topologies

Parameter name	Value range
$K$	5, 15, ..., 75
<i>connectivity</i>	0.1, 0.2, ..., 0.8
$g_k$	50, 250, 450, 650, 850
$bw(l_k)$	10, 30, ..., 90
max-connect	5, 15, ..., 45
$s_k$	100
<i>heterogeneity</i>	0.4, 0.6, 0.8

above parameters. In total, we performed experiments over 29298 generated platform configurations.

**(2) Random graph topologies.** We randomly generated graphs with  $k = 5, 15, \dots, 95$  nodes. Any two nodes are connected with probability *connectivity*, which we vary. As for the Tiers-generated topologies, we assume shortest path routing. We also define a *heterogeneity* parameter that specifies the maximum spread of values for a single parameter across nodes/links in the platform, used as follows. Table 2 shows the mean parameter values used for instantiating the platform configuration. For each combination of these mean values we generate 10 random platform configurations by sampling each platform parameter uniformly between  $mean * (1 - heterogeneity)$  and  $mean * (1 + heterogeneity)$ , where *mean* denotes the corresponding mean value. Note that since only relative values are meaningful in this setting, we fixed the computing speed at  $s_k = 100$ . In total we performed experiments over 269835 random platform configurations.

## 6.2 Results

**6.2.1 LPR** Our first (expected) observation from our simulation results is that LPR always performs poorly, for both Tiers-generated and random graph topologies. In most cases, LPR leaves a significant portion of the network capacity unutilized, and in some cases all  $\beta_{k,l}$  values are actually rounded down to 0, leading to an objective value of 0.

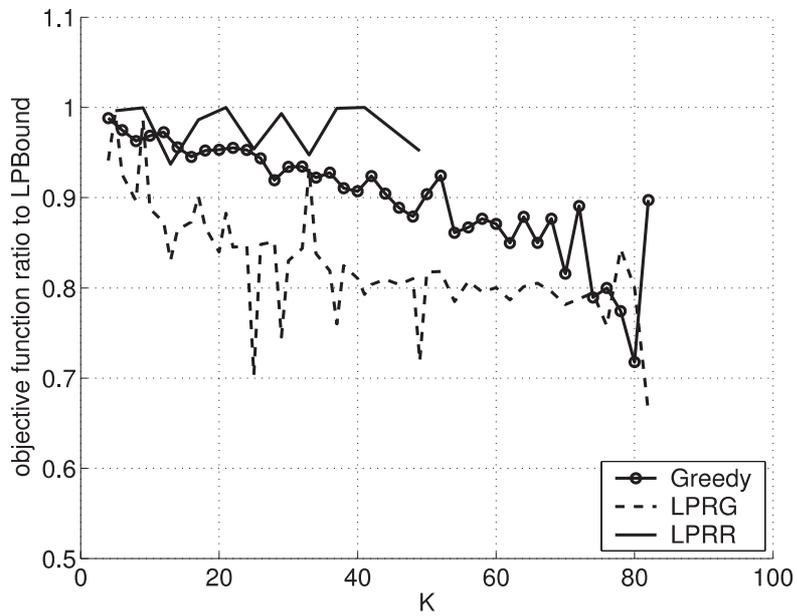
**6.2.2 G vs. LPRG** More interesting is the comparison between G and LPRG. Their relative performance is clearly dependent on topologies.

**(1) Tiers-generated topologies.** In these cases, G performed consistently better than LPRG. Over all platform

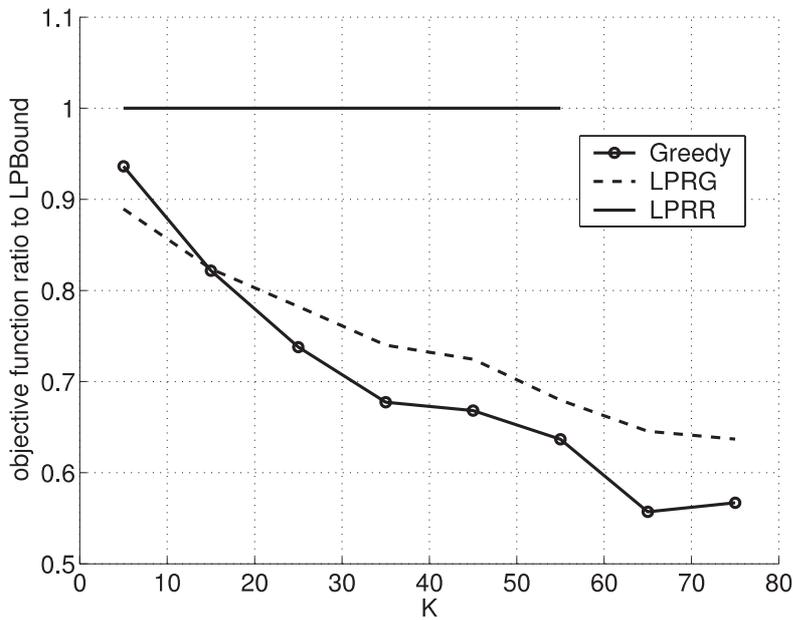
configurations, the average ratio of the objective function achieved by G to that by LPRG was 1.18, with a standard deviation of 31.5, and G was better than LPRG in 81% of the cases. For a closer look, Figure 4(a) plots the average ratio of the objective values achieved by G and LPRG to the upper bound on the optimal obtained by solving the rational linear program, versus the number of clusters  $K$ . We see that G achieved objective values about 5% ~ 10% higher than LPRG in most cases. But as  $K$  increases, both heuristics failed to achieve objective values close to the upper bound on the optimal.

To explain the poor performance of LPRG relative to G, we examined the simulation logs closely. It turns out that, after solving the rational linear program, there are often some clusters that send a portion of their load to other clusters using a rational number of network connections that is strictly lower than 1. After rounding this value down to 0, such clusters then have no opportunity to send off this load portion and become oversubscribed: they take the objective value of the rational linear program down. Conversely, the clusters that were supposed to receive this load are now undersubscribed and have cycles to spare. During the greedy step, the G heuristic sometimes picks an undersubscribed cluster first, which causes this cluster to use up its own spare cycles for its own load. This does not help the MAX-MIN objective value as this cluster was typically better off than the oversubscribed clusters. Furthermore, one or more oversubscribed clusters have now lost the opportunity to use these cycles, which harms the MAX-MIN objective value. In contrast, when the G heuristic starts from scratch, it balances the load better by allowing such undersubscribed clusters to use a full network connection early on in the resource allocation process.

**(2) Random graph topologies.** Over these topologies, the ratio of objective function achieved by LPRG to that by G was 1.98 on average, which is quite different from the results we obtained for Tiers-generated topologies. Figure 4(b) shows that, when there are few nodes in the network, G was slightly better; but as  $K$  increased, LPRG increasingly performed better than G. While these random topologies are not representative of actual networks, it would be interesting to understand which properties of the interconnection topology affect the relative performance of G and LPRG. For now we conjecture that with a less structured topology, such as high-connectivity random graphs, the G heuristic has more opportunity to make bad choices when compared with LPRG, while with a more structured topology, such as tree-structured ones generated by Tiers, the number of choices is more limited and a greedy approach is effective. We did not find any significant trend in the relative performance of G and LPRG with respect to the *heterogeneity* parameter.



(a) Tiers-generated topologies



(b) Random graph topologies

Fig. 4 Performance of G, LPRG and LPRR, relative to the upper bound of the optimal.

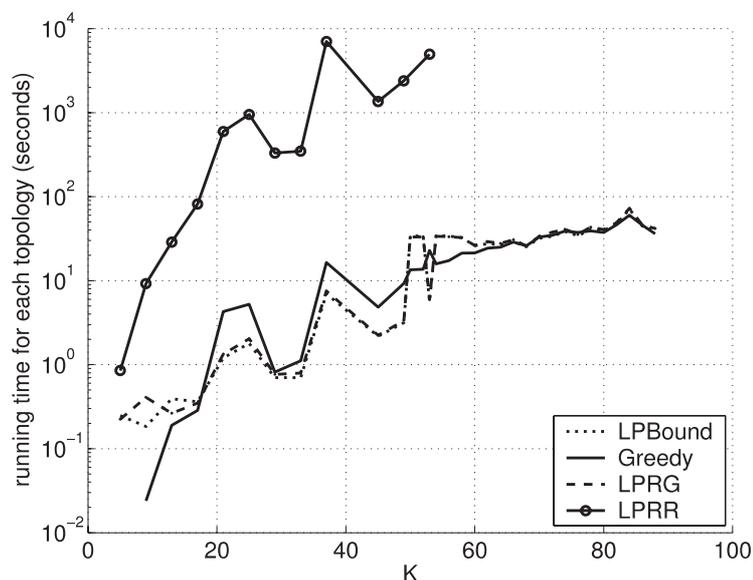


Fig. 5 Running time of G, LPRG and LPRR.

**6.2.3 LPRR** Figure 4(a) and 4(b) show that LPRR performed consistently better than both G and LPRG, for both types of topologies. For Tiers-generated topologies, LPRR achieved an objective value very close to the upper bound of the optimal, even at  $K = 45$ . LPRR almost always achieved the upper bound for the random graph topologies. Explaining why LPRR performed better on random graph topologies than on Tiers-generated topologies is left for future work. Note that since LPRR is much more time consuming than the other heuristics, solving  $K^2$  linear programs, we evaluated LPRR only on a small subset of our topologies.

**6.2.4 Running Time** Figure 5 shows how time-consuming each of our heuristics were and plots their running time in seconds on a 1GHz Pentium processor versus the number of clusters, on a logarithmic scale. These running times were obtained when running our heuristics on the Tiers-generated topologies. We see that LPRR was very expensive: at  $K = 50$ , each run of LPRR took approximately 1 hour, as opposed to 10 to 30 seconds for G and LPRG. This implies that for a real platform with many clusters G and LPRG may be more practical than LPRR.

## 7 Perspectives on Implementation

In this section we discuss how our work and results could be implemented as part of a framework for deploying divisible load applications. Consider a Virtual Organization (VO; Foster, Kesselman, and Tuecke 2001) in which participating sites hold resources that they are willing to contribute for the execution of divisible load applications, as well as users who wish to execute such applications. The G heuristic could be implemented as part of a centralized broker that would manage divisible load applications and the resources they can use, for the entire VO. VO participants would register their resources to the broker, and application requests would be submitted to the broker by users. Note that because our work aims at optimizing steady-state throughput, it provides very good schedules for situations in which applications run for a significant amount of time so that the start-up and clean-up phases of application executions are negligible when compared with the entire application execution time. This is a likely scenario in a VO that supports VO-wide application executions.

The broker needs to gather all relevant information to instantiate the LP formulation of the scheduling problem

(which is needed to implement the G heuristic as well), as given in Section 3.1. Most important are the  $\pi_k$  coefficients that are used to define the objective function. These coefficients define the policies that govern resource sharing over the entire grid, and these policies should be configured at the broker by a VO administrator (or by any kind of contracting system that is in place in the VO). As mentioned in Section 3.1, the objective function can be extended so that different weights are associated for each pair of sites,  $\pi_{k,b}$ , thereby quantifying peering relationships between VO participants. Furthermore, other constraints can be added to the linear program to reflect other arbitrary resource sharing policies (e.g. no more than 10% of resources at cluster  $k$  can be used for applications originating from clusters  $i$  and  $j$ ). The main point here is that our linear program can be refined to express a wide variety of resource sharing policies and peering relationships among VO participants. It would be interesting to see how the G heuristic compares with LPRG and LPRR for more constrained scheduling problems. The broker also needs to be configured so that the number of network connections used by divisible load applications in the VO does not exceed the max-connect threshold. While in this paper we have looked at a general model in which every link has its own threshold, in practice the VO administrator may not have sufficient knowledge of the network topology. In this case VO administrators would just configure a limit on the number of connections on each path between each pair of participating sites, and the same limit could be enforced for each path. The broker needs to have estimates of the compute and transfer speeds that are achieved on the resources. This can be done by querying grid information services (Wolski, Spring, and Hayes 1999; Czajkowski et al. 2001; Ganglia), or by directly observing the performance being delivered by the resources. The latter method may prove easier to implement if divisible load applications are continuously running. The best solution is probably to use both methods and combine them to obtain estimates of achievable performance, as done for instance in (Casanova and Berman 2003). The broker needs to adapt its scheduling decisions as resource availability fluctuates and as applications start and complete. One simple option is to allow adaptation to occur after each scheduling period. Additionally, the schedule could be recomputed on-the-fly as soon as a new application is submitted to the broker or a running application completes.

An intriguing question is that of a decentralized implementation of the broker. In most of today's VOs a centralized broker could probably be engineered in a way that provides appropriate scalability and performance. But in larger VOs the broker could become a performance bottleneck (not to mention it being a single point of failure). In distributed brokering multiple independent brokers would either cooperate, perhaps in a peer-to-peer fashion,

or make their own decisions autonomously. Distributed scheduling is a notoriously difficult question. However, a few simple and elegant solutions have been provided in some specific cases (e.g. see Carter et al. 2003), and it would be interesting to investigate whether a distributed algorithm exists that generates a schedule approaching the solution of our linear program.

Finally, it would be straightforward to provide client-side software by which users could specify, instantiate, and submit their applications to the broker. In fact, the APST software has recently been extended to support divisible load applications (Raadt, Yang, and Casanova 2004). APST provides its own broker that schedules a single application over resources and that handles all deployment logistics for computation, data, resource, and security management. As such, the APST software would provide a good basis on which to build a more general broker that supports multiple applications within a VO.

## 8 Conclusion

We have addressed the steady-state scheduling problem for multiple concurrent divisible applications running on platforms that span multiple clusters distributed over wide-area networks. This is an important problem because divisible load applications are common and make up a significant portion of the mix of grid applications. Only a few authors had previously explored the simultaneous scheduling of multiple such applications on a distributed computing platform (Bharadwaj and Barlas 2002; Wong et al. 2003) and in this paper we have made the following contributions. We defined a realistic platform model that captures some of the fundamental network properties of grid platforms. We then formulated our scheduling problem as a mixed integer-rational linear program that enforces a notion of weighted priorities and fairness for resource sharing between applications. We proposed a greedy heuristic, G, and three heuristics based on the rational solution to the linear program: LPR, LPRG, and LPRR. We evaluated these heuristics with extensive simulation experiments for many random platform configurations whose network topologies were generated by Tiers (Calvert, Doar, and Zegura 1997) or were based on purely random graphs. We found that the G heuristic performs better than LPRG on average for the Tiers-generated topologies, and that its performance relative to an upper bound of the optimal decreases with the number of clusters in the platform. We found that for random graph topologies LPRG outperforms G as soon as the number of clusters becomes larger than 15. Here also, the performance of LPRG decreases as the number of clusters increases. We also found that the LPRR heuristic leads to better schedules than G but at the cost of a much higher complexity, which may make it impractical for large numbers of clusters.

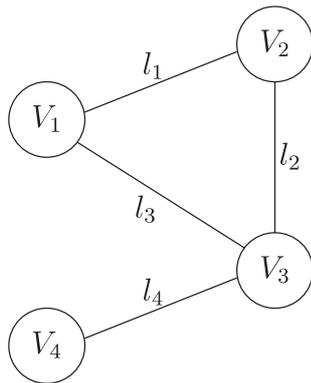
We will extend this work in several directions. First, we will simulate platforms and application parameters that are measured from real-world testbeds and applications suites (Beynon et al. 2002; Spencer et al. 2002). We have gathered such information as part of other research projects. While this paper provides convincing evidence of the relative merit of our different approaches, simulations instantiated specifically with real-world data will provide a quantitative measure of absolute performance levels that can be expected with the best heuristics. Second, we will strive to use an even more realistic network model, which would include link latencies, TCP bandwidth sharing behaviors according to round-trip times, and more precise backbone characteristics. Some of our recent work (see Legrand, Marchal, and Casanova 2003; Casanova 2004) provides the foundation for refining our network model, based both on empirical measurements and on theoretical modeling of network traffic. Finally, one could envision extending our application model to address the situation in which each divisible load application consists of a set of tasks linked by dependencies. This would be an attractive extension of the mixed task and data parallelism approach (Chakrabarti, Demmel, and Yelick 1997; Subhlok and Vondran 2000; Braun et al. 2001) to heterogeneous clusters and grids.

## Appendix

We give the proof of Theorem 1 in this section.

**Theorem 1.** *STEADY-STATE-DIVISIBLE-LOAD*( $\mathcal{P}$ ,  $\pi$ ,  $\rho$ ) is NP-complete.

**Proof.** We first prove that this problem belongs to NP. Given an instance  $\mathcal{I}$  of STEADY-STATE-DIVISIBLE-LOAD, we verify that  $\mathcal{A} = (\alpha, \beta)$  is a valid allocation



**Fig. 6** Example of instance  $\mathcal{I}_1$  of MAXIMUM-INDEPENDENT-SET.

by checking that Equations 6 are satisfied, and that  $\min_k \left\{ \frac{\alpha_k}{\pi_k} \right\} \geq \rho$ , which can be done in polynomial time.

To prove the completeness of STEADY-STATE-DIVISIBLE-LOAD, we proceed by a reduction from MAXIMUM-INDEPENDENT-SET, which is known to be NP-complete (Garey and Johnson 1979). Consider an arbitrary instance  $\mathcal{I}_1$  of MAXIMUM-INDEPENDENT-SET: given a non-oriented graph  $G = (V, E)$  and an integer bound  $B$ , does there exist a subset  $V'$  of  $V$  of cardinal at least  $B$  and such that no two vertices of  $V'$  are joined by an edge of  $E$ ? From  $\mathcal{I}_1$ , we construct the following instance  $\mathcal{I}_2$  of STEADY-STATE-DIVISIBLE-LOAD:

- Let  $V = \{V_1, \dots, V_n\}$ . The platform of  $\mathcal{I}_2$  consists of  $n + 1$  clusters  $C^0, C^1, \dots, C^n$ . A set  $Route(i)$  is associated with each cluster  $C^i$  and will be used to determine the routing in the platform as explained below.
- $E = \{e_1, \dots, e_m\}$ . For each edge  $e_k = (V_i, V_j) \in E$ , we add in the platform graph:
  - two routers  $Q_k^a$  and  $Q_k^b$ ,
  - a backbone link between them:  $l_k^{common} = (Q_k^a, Q_k^b)$ , with  $max-connect(l_k) = 1$  and  $bw(l_k) = 1$ ,
  - a new element  $k$  in  $Route(i)$  and  $Route(j)$ .
 Then, for each set  $Route(i) = \{k_1, \dots, k_{|Route(i)|}\}$  associated to cluster  $C^i$ , we add the following backbone links, all with  $max-connect(l) = 1$  and  $bw(l) = 1$ :

$$l_1^i = (C^0, Q_{k_1}^a),$$

$$l_j^i = (Q_{k_j}^b, Q_{k_{j+1}}^a), \text{ for } j = 1, \dots, |Route(i)|,$$

$$l_{|Route(i)|+1}^i = (Q_{k_{|Route(i)|}}^b, C^i)$$

- Finally, the routing between cluster  $C^0$  and cluster  $C^i$  is given by the following routing path:

$$L_{0,i} = \{l_1^i, l_{k_1}^{common}, l_2^i, l_{k_2}^{common}, \dots, l_{k_{|Route(i)|}}^{common}, l_{|Route(i)|+1}^i\} \quad (7)$$

- Cluster  $C^0$  has specific characteristics:  $g_0 = n$  and  $s_0 = 0$ , while all other clusters are such that  $g_i = s_i = 1$ .
- We let  $\delta_k = w_k = 1$  for each application  $A_k$ , and we set the priority factors to  $\pi_0 = 1$  and  $\pi_i = 0$  for  $j = 1, \dots, n$  ( $C^0$  is the only cluster which has work to do).
- The throughput bound  $\rho$  of  $\mathcal{I}_2$  is set to  $B$ .

The platform graph that we have constructed has a strong property which is expressed in the following lemma:

**Lemma 1.** Two routes  $(C^0, C^i)$  and  $(C^0, C^j)$  in the platform graph of instance  $\mathcal{I}_2$  share a common backbone link if and only if the edge  $(V_i, V_j)$  belongs to the graph  $G$  of instance  $\mathcal{I}_1$ .

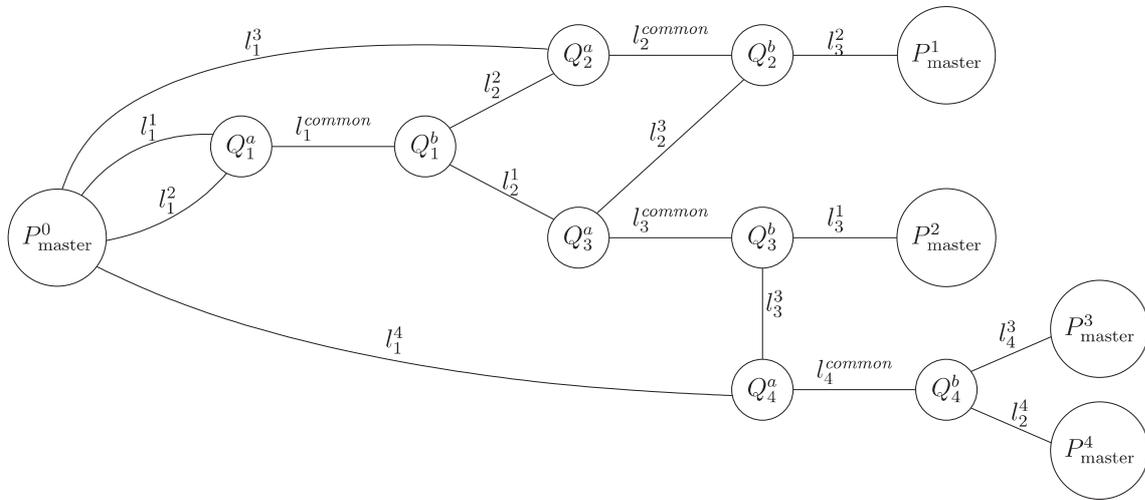


Fig. 7 Instance  $\mathcal{I}_3$  of STEADY-STATE-DIVISIBLE-LOAD built from the previous instance example  $\mathcal{I}_1$ .

**Proof.** Assume first that edge  $e_k = (V_i, V_j)$  belongs to  $G$ . Then, by construction, as  $k$  is added to the list  $Route(i)$  and  $Route(j)$ , the corresponding link  $l_k^{common}$  belongs both to  $L_{0,i}$  and  $L_{0,j}$ ; the routes  $(C^0, C^i)$  and  $(C^0, C^j)$  share the common link  $l_k^{common}$ .

Assume now that routes  $(C^0, C^i)$  and  $(C^0, C^j)$  share a backbone link. According to Equation 7, this link is a  $l_k^{common}$  for some  $k$ . As  $l_k^{common} \in L_{0,i}$  and  $l_k^{common} \in L_{0,j}$ , then  $k \in Route(i)$  and  $k \in Route(j)$ . The construction of these sets shows that there is an edge  $e_k$  between  $V_i$  and  $V_j$  in  $G$ .

We now prove that there exists a solution to  $\mathcal{I}_1$  if and only if there exists a solution  $\mathcal{I}_2$ :

- Assume that there exists an independent set  $V'$  solution of  $\mathcal{I}_1$  (so  $|V'| \geq B$ ). From  $V'$ , we construct the following allocation  $\mathcal{A}$ :

$$\forall i, \alpha_{0,i} = \beta_{0,i} = \begin{cases} 1 & \text{if } V_i \in V' \\ 0 & \text{otherwise} \end{cases}$$

$$\forall j \neq 0, \forall i, \alpha_{j,i} = \beta_{j,i} = 0$$

As  $V'$  is an independent set, there is no edge in  $G$  between any two vertices of  $V'$ , so there is no common backbone link between the routes defined by non-zero

values of the  $\beta$ 's. Each backbone is used by at most one route, and since  $max-connect = 1$  for all backbones, Equation 3 is satisfied. There are  $|V'|$  (which is less than  $n$ ) different routes outgoing from  $C^0$ , and none incoming to it, so Equation 2 is fulfilled since  $g_0 = n$ . For all other clusters  $C^i$  ( $i > 0$ ), at most one route with bandwidth 1 is incoming and none is outgoing, so Equation 2 is satisfied since  $g_i = 1$ . Each cluster  $C^i$  such that  $V_i \in V'$  has to compute an amount of work of 1 unit, which is not more than its speed, so Equation 1 is satisfied.

Hence,  $(\alpha, \beta)$  defines a valid allocation which reaches the throughput of  $|V'| \geq B = \rho$ . This is a solution for  $\mathcal{I}_2$ .

- Assume now that  $(\alpha, \beta)$  is a solution of  $\mathcal{I}_2$ , which means that it is a valid allocation whose throughput is at least  $\rho$ . As  $C^0$  has no computing power, it has to delegate the work to other clusters. Each other cluster has a computing speed of one task every time-unit, so there exist at least  $\rho$  different routes from cluster  $C^0$  to  $\rho$  distinct clusters  $C^{k_1}, \dots, C^{k_\rho}$ . Since  $max-connect = 1$  for each backbone link, only one route can go through each backbone link. Hence, for every couple of routes to clusters  $C^{k_i}$  and  $C^{k_j}$ , no link is shared, which means that there is no edge  $(V_i, V_j)$  in the original graph. So the set of the corresponding vertices  $V' = \{V_{k_1}, \dots, V_{k_\rho}\}$  is an independent set in  $G$ . As the cardinal of this set is  $\rho = B$ ,  $V'$  is a solution of the instance  $\mathcal{I}_1$ .

## Acknowledgements

This paper is based upon work supported by the National Science Foundation under Grant No. 0234233.

## Author Biographies

*Loris Marchal* is currently a PhD student in the LIP laboratory at ENS Lyon. He is mainly interested in parallel algorithm design for heterogeneous platforms and in scheduling techniques.

*Dr. Yang Yang* currently works at Amazon.com. His research interests are scheduling of parallel and distributed applications. He obtained his BE from Tsinghua University, PR China in 2000, and his MS and PhD from the University of California, San Diego in 2003 and 2005.

*Dr. Henri Casanova* is an Assistant Professor in the Information and Computer Science Department at the University of Hawaii at Manoa. His research interests are in the area of parallel, distributed, grid, and internet computing. In particular, his research emphasizes the modeling and the simulation of emerging platforms and applications, as well as both the theoretical and practical aspects of scheduling problems. He obtained his BS from the Ecole Nationale Supérieure d'Electronique, d'Electrotechnique, d'Informatique et d'Hydraulique de Toulouse, France in 1993, his MS from the Université Paul Sabatier, Toulouse, France in 1994, and his PhD from the University of Tennessee, Knoxville in 1998. Prior to joining the University of Hawaii, he was a Research Scientist at the San Diego Supercomputer Center and an Adjunct Professor in the Department of Computer Science and Engineering at the University of California, San Diego.

*Prof. Yves Robert* received the PhD degree from Institut National Polytechnique de Grenoble in 1986. He is currently a full professor in the Computer Science Laboratory LIP at ENS Lyon. He is the author of four books, 95 papers published in international journals, and 120 papers published in international conferences. His main research interests are scheduling techniques and parallel algorithms for clusters and grids. He is a Fellow of IEEE.

## References

- Allcock, W., Bester, J., Bresnahan, J., Chervenak, A., Liming, L., and Tuecke, S. 2001. GridFTP: protocol extension to FTP for the grid. Grid Forum Internet-Draft.
- Ausiello, G., Crescenzi, P., Gambosi, G., Kann, V., Marchetti-Spaccamela, A., and Protasi, M. 1999. *Complexity and Approximation*. Berlin, Germany: Springer.
- Banino, C., Beaumont, O., Carter, L., Ferrante, J., Legrand, A., and Robert, Y. 2004. Scheduling strategies for master-slave tasking on heterogeneous processor platforms. *IEEE Transactions on Parallel Distributed Systems* 15(4):319–330.
- Barlas, G. 1998. Collection-aware optimum sequencing of operations and closed-form solutions for the distribution of a divisible load on arbitrary processor trees. *IEEE Transactions on Parallel Distributed Systems* 9(5):429–441.
- Bataineh, S., Hsiung, T., and Robertazzi, T. G. 1994. Closed form solutions for bus and tree networks of processors load sharing a divisible job. *IEEE Transactions on Computers* 43(10):1184–1196.
- Beaumont, O., Casanova, H., Legrand, A., Robert, Y., and Yang, Y. 2003. Scheduling divisible loads for star and tree networks: main results and open problems. Technical Report RR-2003-41, LIP, ENS Lyon, France. To appear in *IEEE Transactions on Parallel and Distributed Systems*.
- Beaumont, O., Legrand, A., and Robert, Y. 2003. Scheduling divisible workloads on heterogeneous platforms. *Parallel Computing* 29:1121–1152.
- Bertsekas, D. and Gallager, R. 1987. *Data Networks*. London: Prentice Hall.
- Bertsimas, D. and Gamarnik, D. 1999. Asymptotically optimal algorithm for job shop scheduling and packet routing. *Journal of Algorithms* 33(2):296–318.
- Beynon, M. D., Kurc, T., Sussman, A., and Saltz, J. 2002. Optimizing execution of component-based applications using group instances. *Future Generation Computer Systems* 18(4): 435–448.
- Bharadwaj, V. and Barlas, G. 2002. Efficient scheduling strategies for processing multiple divisible loads on bus networks. *Journal of Parallel and Distributed Computing* 62: 132–151.
- Bharadwaj, V., Ghose, D., Mani, V., and Robertazzi, T. 1996. *Scheduling Divisible Loads in Parallel and Distributed Systems*. Los Alamitos, CA: IEEE Computer Society Press.
- Bharadwaj, V., Ghose, D., and Robertazzi, T. 2003. Divisible load theory: a new paradigm for load scheduling in distributed systems. *Cluster Computing* 6(1):7–17.
- Braun, T., Siegel, H., Beck, N., Bölöni, L., Maheswaran, M., Reuther, A., Robertson, J. P., Theys, M., Yao, B., Henggen, D., and Freund, R. F. 2001. A comparison of eleven static heuristics for mapping a class of independent tasks onto heterogeneous distributed computing systems. *Journal of Parallel and Distributed Computing* 61(6):810–837.
- CAF. CDF Analysis Farms. <http://cdfcaf.fnal.gov/>.
- Calvert, K., Doar, M., and Zegura, E. W. 1997. Modeling internet topology. *IEEE Communications Magazine* 35:160–163.
- Carter, L., Casanova, H., Ferrante, J., and Kreaseck, B. 2003. Autonomous protocols for bandwidth-centric scheduling of independent-task applications. *International Parallel and Distributed Processing Symposium IPDPS'2003*. Los Alamitos, CA: IEEE Computer Society Press.
- Casanova, H. 2004. Modeling Large-Scale Platforms for the Analysis and the Simulation of Scheduling Strategies. *Proceedings of the 6th Workshop on Advances in Parallel and Distributed Computational Models (APDCM)*.
- Casanova, H. and Berman, F. 2003. Parameter Sweeps on the Grid with APST. *Grid Computing: Making The Global*

- Infrastructure a Reality*, A. Hey, F. Berman, and G. Fox, (Eds). Chichester: John Wiley.
- Chakrabarti, S., Demmel, J., and Yelick, K. 1997. Models and scheduling algorithms for mixed data and task parallel programs. *Journal of Parallel and Distributed Computing*, 47:168–184.
- Coudert, D. and Rivano, H. 2002. Lightpath assignment for multifibers WDM optical networks with wavelength translators. *IEEE Global Telecommunications Conference (Globecom'02)*, Session OPNT-01-5. Los Alamitos, CA: IEEE Computer Society Press.
- Czajkowski, K., Fitzgerald, S., Foster, I., and Kesselman, C. 2001. Grid information services for distributed resource sharing. *Proceedings of the 10th IEEE Symposium on High-Performance Distributed Computing (HPDC-10)*.
- Foster, I., Kesselman, C., and Tuecke, S. 2001. The anatomy of the grid: enabling scalable virtual organizations. *International Journal of High Performance Computing Applications* 15(3).
- Ganglia. The Ganglia Project. <http://ganglia.sourceforge.net>.
- Garey, M. R. and Johnson, D. S. 1979. *Computers and Intractability, a Guide to the Theory of NP-Completeness*. W. H. Freeman and Company.
- Ghose, D. and Robertazzi, T. (Eds) 2003. *Special issue on Divisible Load Scheduling*. *Cluster Computing* 6(1).
- Lee, C. and Stepanek, J. 2001. On Future Global Grid Communications Performance. *HCW'2001, the 10th Heterogeneous Computing Workshop*. Los Alamitos, CA: IEEE Computer Society Press.
- Legrand, A., Marchal, L., and Casanova, H. 2003. Scheduling Distributed Applications: The SimGrid Simulation Framework. *Proceedings of the Third IEEE International Symposium on Cluster Computing and the Grid (CCGrid'03)*.
- Motwani, R., Naor, J., and Raghavan, P. 1996. *Approximation Algorithms for NP-Hard Problems*, chapter 11. Boston, MA: PWS Publishing Company.
- Raadt, K. van der, Yang, Y., and Casanova, H. 2004. APST-DV: Divisible Load Scheduling and Deployment on the Grid. Technical Report CS2004-0785, Dept. of Computer Science and Engineering, University of California, San Diego.
- Robertazzi, T. (Undated) Divisible Load Scheduling. <http://www.ece.sunysb.edu/~tom/dlt.html>.
- Robertazzi, T. 1993. Processor equivalence for a linear daisy chain of load sharing processors. *IEEE Transactions on Aerospace and Electronic Systems* 29:1216–1221.
- Robertazzi, T. 2003. Ten reasons to use divisible load theory. *IEEE Computer* 36(5):63–68.
- Spencer, M., Ferreira, R., Beynon, M., Kurc, T., Catalyurek, U., Sussman, A., and Saltz, J. 2002. Executing multiple pipelined data analysis operations in the grid. *2002 ACM/IEEE Supercomputing Conference*. ACM Press.
- Subhlok, J. and Vondran, G. 2000. Optimal use of mixed task and data parallelism for pipelined computations. *Journal of Parallel and Distributed Computing*, 60:297–319.
- Wolski, R., Spring, N., and Hayes, J. 1999. The network weather service: a distributed resource performance forecasting service for metacomputing. *Future Generation Computer Systems* 15(10):757–768.
- Wong, H., Yu, D., Bharadwaj, V., and Robertazzi, T. 2003. Data intensive grid scheduling: multiple sources with capacity constraints. *PDCS'2003, 15th International Conference on Parallel and Distributed Computing and Systems*. IASTED Press.
- Yu, D. and Robertazzi, T. 2003. Divisible load scheduling for grid computing. In *PDCS'2003, 15th International Conference on Parallel and Distributed Computing and Systems*. IASTED Press.