

# Hierarchical Scheduling in Grid Systems

Khaldoon AlZoubi  
Carleton University, Ottawa, Canada  
Email: [Kalzoubi0708@rogers.com](mailto:Kalzoubi0708@rogers.com)

Sivarama Dandamudi  
Carleton University, Ottawa, Canada  
Email: [sivarama@scs.carleton.ca](mailto:sivarama@scs.carleton.ca)

## ABSTRACT

This paper proposes hierarchal scheduling schemes for grid systems: a self-discovery scheme for the resource discovery stage and an adaptive child scheduling method for the resource selection stage. In addition, we propose three rescheduling algorithms: the butterfly, fallback and load-balance. We also propose a hybrid system to combine the proposed hierarchal schemes with the well-known peer-to-peer (P2P) principle. We compare the performance of the proposed schemes against the P2P systems with respect to a set of predefined metrics.

Key Words: Grid Computing, Grid systems, Hierarchal Scheduling, P2P systems, Cluster Computing, Parallel Scheduling, High-Performance Computing.

## 1. Introduction

The current status of computation is equivalent in some respects to the status of the electricity in around 1910s [4]. At that time, electrical power was generated by generators for specific individuals or organizations needs. Truly, the real electricity influence in our lives was born with the creation of the electric power grid which was provided via sharing generators. The “Grid Computing” term was adopted from the electricity grid to *amplify computational power* via sharing computational resources since both grids are similar with respect to their infrastructure and purpose. The term “the Grid” started in the mid-1990s [4, 5] to portray infrastructure for both scientific and commercial distributed computing communities and has been gaining popularity ever since. The “Grid” can be defined as a parallel and distributed system that enables large collection of geographically distributed heterogeneous systems that usually span over several organizations to share a variety of resources dynamically at runtime depending on their availability, capability,

user’s requirements and any other predefined rules by local systems and resources owners. The type of sharing in the grid gives the impression of a powerful self-managing virtual computer. The Internet can be an ideal choice to link thousands or millions of computers since it already connects the whole world – if a node’s IP address is known, it can then receive data from another node. Benefits of grids can be extensive. They include: (1) expanding computing power, since grids unleash the hidden computing power that is not being used for most of the time (e.g. most machines in a typical organization are busy less than 5% of the time [2]), (2) improving productivity and collaboration among organizations (i.e. wider audience) in a dynamic and geographically distributed manner to form one powerful computing system, and (3) solving complex problems that were previously unsolvable.

The rest of the paper is organized as follows. In Section 2, the grid scheduling stages and some of the grid challenges are described. Section 3 presents the self-discovery method to be used in the resource discovery stage and the adaptive hierarchical scheduling (AHS) method to

be used in scheduling jobs on selected resources. Note that the AHS method is based on the AHS method for parallel and cluster systems presented in [3]. In addition, we present three rescheduling dynamic algorithms: the butterfly, the fallback and the load-balance algorithms. In Section 4, the simulation model and samples of the preliminary results are given. Conclusions are given in Section 5. Refer to [1] for more in depth discussion of the presented schemes and the complete set of results.

## 2. Grid Scheduling Stages

Grid characteristics must be taken into account to be able to perform efficient scheduling. Grid schedulers must make scheduling decisions in a very challenging environment where it has: (1) no control over the resources since they don't own them, (2) distributed resources, (3) dynamic existence of resources (i.e. resources may be added or removed from the grid at any time), (4) dynamic information collection, (5) heterogeneous resources (jobs must match appropriate resources in order to be executed as requested by the users), (6) tentative scheduling until the allocation of actual resources (i.e. resources may be taken from the grid before a job actually uses them).

In general, grid scheduling is performed in three stages [7]. (1) Resource discovery stage, which produces a set of matched resources. Schedulers are expected to collect static information (e.g. operating systems), from local schedulers or general information systems (GIS), in order to be able to perform job matching. (2) Resources are selected (i.e. resource selection stage), from the list obtained during the first stage, and are expected to meet user's imposed constraints (e.g. deadlines). Schedulers are expected to collect dynamic information (e.g. system load) for this stage (3) Transfer jobs to selected resources (i.e. job execution stage).

## 3. Hierarchical Grid Scheduling

Grid schedulers are structured in a tree form that we call grid tree, as shown in Figure 3.1, where grid schedulers (GS) are placed into the tree according to their geographical locations. Users submit their jobs, in the form of requests to the grid via the grid system scheduler (GSS), which is the root node of the grid tree. A user's request describes the job in terms of the job minimum requirements (JMR) (e.g. operating system) in order to be matched to resources and any other constraints imposed by the user (e.g. completion deadline). A leaf grid scheduler (LGS), which is a node on top of local scheduler(s), connects directly with the user's workstation, brings the physical job to the grid (once a job is about to be mapped to the selected resources) and serves as a middleware between the user's workstation and the allocated resources. Note that LGSs may be combined with local schedulers in one unit.

Theoretically, a scheduler, in our proposed systems, can break grid jobs into several subjobs to be executed in parallel on different children's partitions. However, the art of automatic transformation into parallelism of an arbitrary job is in its infancy stage [2]. In our case, we assume whole jobs are submitted to resources.

### 3.1 Resource Discovery Stage

In this stage, we propose the *self-discovery* method. The purpose of this method is to produce a set of logical channels to be used as paths by jobs (in the next scheduling stage) to get to their physical resources. Logical channels serve as a map for jobs so that they know how to reach resources that can meet their computational requirements. The self-discovery method omits irrelevant dissimilarities between resources of different sites. The principle behind this method is that resources are equivalent to each other if they match the same set of jobs. For example, one site advertises INTEL architectures and another site advertises AMD architectures. Now, suppose the grid has a set of jobs that can be executed on either INTEL or AMD

platform. In this case, the grid system considers architecture INTEL as equivalent to architecture AMD for those jobs in the set since they can be executed on either platform. However, suppose now another set of jobs only requires architecture AMD in order to execute. In this case, the grid system considers architecture INTEL as nonequivalent to architecture AMD for the later set of jobs since those jobs can only be executed on the AMD platform.

LGSs collect and store all static information about resources either directly from local schedulers or from a GIS. Thus, information about local resources is distributed across the grid which leads to (1) increasing system scalability and (2) maintaining up-to-date databases. LGSs also initiate the resource discovery stage at system startup or when no jobs in the grid tree match their advertised resources by issuing the request for job matching (RFJM) message to their parents, which in turn forward the RFJM messages to the grandparents, and so on until the RFJM message is received by the GSS, enabling it to initiate resource discovery to all of its raw jobs (i.e. new jobs that haven't previously been through resource discovery stage). However, if the GSS has no raw jobs, it will then backlog the RFJM message until receiving new jobs.

The GSS starts the resource discovery stage by: (1) broadcasting a special message to all of its children to destroy all channels in the system, and (2) passing all raw jobs to all of its children as one block. The children in turn pass the raw jobs as one block to the grandchildren, and so on until they reach the LGSs at the bottom of the grid tree. LGSs match raw jobs requirement to their local resources and inserts all raw jobs that match resources in their bags. They will then pass a copy of their bags (along with any previous matched jobs) to their parents. Note that (1) LGSs save all requests that they receive from their parents regardless if they have matched or not, enabling LGSs to perform rematching, if needed, due to resources change (i.e. the GSS removes all stored requests from all bags once they are

executed), (2) intermediate schedulers (GS) always pass one RFJM message to their parents on behalf of their children and suppress other RFJMs preventing the GSS from initiating any unnecessary resource discovery, (3) A logical channel is created for every unique jobs bag.

Once a scheduler receives a bag (from one of its children) that is similar to another bag of an existing channel, it will then: (1) create a new branch from its existing channel and bind it with the child's channel, (2) recalculate its channel's processing power based on the new created branch, (3) inform the child with its channel's port number, and (4) update parent, if any, with the new processing power of its channel. However, if the received bag is distinctive, scheduler will then: (1) create a new channel with a new port number, (2) create a new branch from its new channel and bind it with the child's channel, (3) initialize the channel processing power based on the new created branch, (3) inform the child with the channel's port number, and (4) update parent, if any, with the new bag and the channel's port number.

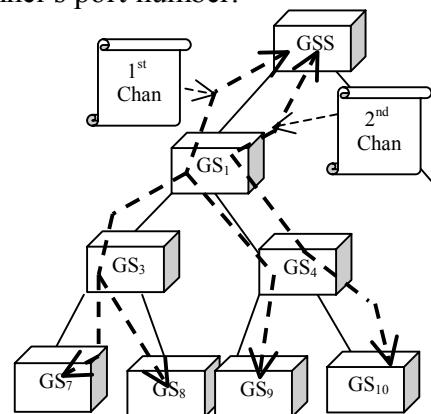


Figure 3.1: A grid tree with two channels

Now consider, as an example, the grid tree shown in Figure 3.1. Suppose the GSS pass a six raw jobs  $J_1$  through  $J_6$  as one block to all of its children. Suppose further that  $J_4$ ,  $J_5$  and  $J_6$  do not match resources at  $GS_{10}$ . In this case, the GSS ends up with two channels (via  $GS_1$ ) where all jobs may use the first channel. On the other hand,  $J_1$ ,  $J_2$  and  $J_3$  are the jobs that can only use the second channel.

### 3.2 Resource Selection Stage

The grid AHS scheduling method, in this stage, uses *self-scheduling* by exploring the parent-child relationship. When a non-root GS wants some work to do, it initiates self-scheduling by sending a request for computation (RFC) message to its parent (via a channel) requesting computation from it. If the parent GS doesn't have computations that can be pushed on that child's channel at the time of receiving the RFC, it in turn generates its own RFC and sends it to its parent on the next level of the grid tree. This process is recursively followed until either the RFC reaches the grid system scheduler (GSS) or a GS with computations is encountered along the path. Note that intermediate schedulers send one RFC message per channel to their parents, but still mark all branches that have received RFCs from. Schedulers use space-sharing policy to distribute computations among channels, upon receiving an RFC message from a channel's branch of a child, as follows:

$$B_{share} = \lceil B_{rate} \times N \rceil$$

where  $B_{share}$  is branch's share of all jobs within scheduler's subtree,  $B_{rate}$  is the branch's transfer rate, and  $N$  is the number of jobs within a scheduler's subtree.  $B_{rate}$  is calculated as follows:

$$B_{rate} = \frac{B_{pwr}}{\sum_{i=1}^M (C_{pwr})_i}$$

where  $B_{pwr}$  is the branch's processing power,  $C_{pwr}$  is the channel's processing power (i.e. total processing power for all of its branches), and  $M$  is the number of channels in a scheduler. Note that, in our case, the channel's processing power is the number of CPUs that resides under that channel, since we assume that our computational resources are parallel computers (see Section 4). However, in reality, we expect processing power to consider more factors such as RAMs, bandwidth, etc.

Now, once a scheduler determines the number of jobs that will be pushed onto a channel's branch, it builds a list of those jobs as one block and pushes them onto that branch. Suppose, as an example, that  $GS_1$ , in Figure 3.1, has 9 jobs upon receiving an RFC message from the first branch of the first channel (i.e. via  $GS_3$ ). Suppose further that the three branches that connect  $GS_1$  with its two children have equivalent processing power. In this case,  $GS_1$  may then push up to 3 jobs onto that branch.

Schedulers perform the following steps to collect the jobs (in order to be pushed onto a channel's branch): (1) invokes the "butterfly" algorithm, (2) collects jobs from the unassigned (i.e. unpushed) jobs, and (3) invokes the "load-balance" algorithm. Note that we expect schedulers, in practice, to collect more dynamic information related to performance (e.g. load) or economics (e.g. prices).

#### 3.2.1 Butterfly Algorithm

The principle behind this scheme is to reschedule jobs to better resources (with respect to predefined metrics) when they become available. Note that this algorithm can be extended to any soft conditions imposed by a user where soft conditions are the ones that the user is willing to live without them until they become available (or if they ever become available). In our case, we consider the geographical closeness of resources with respect to workstations as our metric (we use IP addresses to determine location closeness). Interestingly, a job may keep jumping (like a butterfly) among children's partitions until it settles on the closest resources. In this algorithm, after a scheduler receives an RFC message from a child, the scheduler will then (1) cancel any assigned jobs from other children (if the new child is closer to those jobs workstations), and (2) push them into the new available child's partition.

#### 3.2.2 Fallback Algorithm

The fallback algorithm is intended to reschedule jobs that become incomputable

because of the grid losing their required resources on their scheduled partitions. When an LGS detects resources change, it performs rematching for all saved requests. Now, if an LGS ends up with the same jobs bag, this resource change is then irrelevant. However, if it produces a different bag, it will then pass it onto its parent.

Schedulers handle received bags in this stage as previously described in the resource discovery stage. Additionally, schedulers mainly have to carry out the following (of course, parents will also be updated): (1) remove any jobs that become incomputable, (2) recalculate modified channels processing power, and (3) delete any broken channels. For example, assume  $GS_{10}$  changes resources and produces similar bag to  $GS_9$ 's bag in Figure 3.1. In this case,  $GS_4$  will then connect  $GS_{10}$  to the first channel and inform its parent ( $GS_1$ ) of two things: (1) the first channel new updated processing power and (2) the broken second channel. Now if  $GS_4$  has jobs that become incomputable, it will then remove them and update  $GS_1$ . Note that  $GS_1$  reschedules those returned jobs (if they still computable on its partition) with a priority (in our case, the lesser the sequence number, the higher the priority). For instance,  $GS_1$  may swap some of those returned jobs with assigned jobs in  $GS_3$ 's partition in order to execute jobs in the same order of their arrival to the grid.

### 3.2.3 Load-Balance Algorithm

As stated earlier, schedulers determine the number of jobs (that will be pushed into a channel's branch) by considering all jobs within their subtrees. They will then collect those jobs by invoking the butterfly algorithm and from queued unassigned jobs. Schedulers will then balance the load among the channels by canceling already assigned jobs and then reschedule them on the channel's branch that just requested more work. Schedulers are always responsible for balancing all assigned jobs among their children's channels within their subtrees, since a parent may cause a

child's subtree to get imbalanced (i.e. of course, parents do not know how assigned jobs are distributed within their children's subtrees). For example, assume that each of  $GS_7$  and  $GS_8$ , in Figure 3.1, has four queued jobs (of course,  $GS_1$  assumes that all 8 jobs are still queued at  $GS_3$ ). Suppose now that  $GS_1$  decides to cancel four jobs to reschedule them on a  $GS_4$ 's channel in order to balance its subtree. Now, if  $GS_1$  cancels the four jobs queued at  $GS_8$ , it imbalances the  $GS_3$ 's subtree. In this case,  $GS_3$  will balance its subtree upon receiving an RFC message from  $GS_8$  (or will forward the RFC message to  $GS_1$ , if the RFC message is received from  $GS_7$ ).

### 3.3 Hybrid System

The hierarchal system (one grid tree) has major drawbacks: (1) all requests are submitted to the GSS, which may become overwhelmed with too many requests, (2) difficult to bring too many organizations to agree on things such as constructing the grid tree, controlling GSS policies (e.g. security), dealing with new joined organizations, and (3) difficult to convince companies to replace their peer-to-peer (P2P) based grid systems. The hybrid system, which is several grid trees that act also as peers to each other, as shown in Figure 3.2, does not only overcome the above drawbacks, but also provides organizations more efficient ways to manage their own resources such as stamping foreign requests with low priority, isolating their resources swiftly from the entire grid without pumping out their pending requests or using their resources, etc.

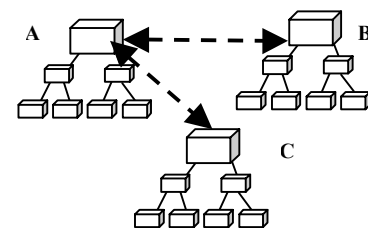


Figure 3.2: Hybrid System Example

In the hybrid system, a GSS in a grid tree also acts as a peer GSS (PGSS) that

forwards requests (after decrementing hop count) to its neighbors. Of course neighbors can be part of any other system types (e.g. P2P system). The P2P system can be viewed as a hybrid system where each grid tree has only one scheduler, and the hierarchal system can be viewed as a hybrid system with one peer. In our case, we assume that (1) requests are always forwarded to neighbors (i.e. a request dies when hop count reaches 0), hence, grid trees also serve as backups to each other, and (2) foreign and home requests are queued in the same fashion.

#### **4. Simulation Model and Results**

This section presents the simulation model and samples of the obtained preliminary results. Readers are encouraged to refer to [1] for the complete set of results and more in depth discussion.

The grid simulation model is broken into three submodels: communication, node and system models, each of which is described below.

##### **4.1 Communication Model**

The communication model, which is used by nodes (i.e. node model) to communicate with each other, consists of 2400 nodes that span across four backbones. Each backbone (600 nodes) consists of four Nets where each Net consists of 10 networks, each network consisting of 15 nodes. Therefore, there are (4 backbones) (4 Nets) (10 networks) (15 nodes) which adds up to 2400 nodes in total in the model. The communication model uses the Discrete Event Simulation (DEVS) CD++ simulator [1, 11] to simulate all communication aspects among all nodes.

The model presented numerous of challenges that we had to address to bring it closer as much as possible to the actual communication over the Internet, which is almost an impossible job to do, since the Internet is a very large unpredictable public network. We've assumed that 10% of the model's capacity accounts for the external Internet load and the routers processing time based on the studies in [1,

8], which were based on actual statistics by the Internet Service Providers (ISP). We've also assumed 64 kilobytes TCP window size to control data flow [1]. Backbones in the model are connected with 1000 km, 9.6 Gb/s (e.g. OC-192 link) cables, Nets/sites are connected with 50 km, 155 Mb/s (e.g. OC3 link) cables [1, 8], and nodes within a site are connected with 100 MBytes/s [1] cables.

##### **4.2 Node Model**

A node, in the communication model, is simply a computer with an IP address. On the other hand, the node component contains the implementation of the proposed schemes for the grid systems in this paper. A node can be configured to operate as a peer, local scheduler, intermediate GS, LGS, GSS, PGSS or a workstation. Note that the node's configuration determines the system model type.

##### **4.3 System Model**

The grid model can be configured to a peep-to-peer (P2P), a hierarchal (one grid tree) or a hybrid (several grid trees) system model, as discussed below in this section. The P2P systems are distributed systems and the only ones, to our knowledge, that are currently well thought-out by researchers [1, 7, 9] to replace the centralized systems. As a result, it is a reasonable choice to be compared against the proposed hierarchal systems in this paper.

###### **4.3.1 Peer-to-Peer (P2P) System**

Once a job request is received at a peer that meets its requirement, it contacts the workstation to offer its service and times for the workstation response (100 ms, in our case). Workstations may accept peers service by responding to it or may refuse peers service by simply not responding to it. If a peer can not accept a workstation request, it decrements the hop count (1000 in our case) in the message and forwards it to its neighbors. In our model, peers accept requests if they meet their deadline, which is three times the estimated execution time

for that job. To improve P2P performance: (1) if a workstation doesn't get a service offer within 2 minutes; it resubmits the job request again to the grid and (2) neighbors are manually configured to be geographically close. However, this may not be the case in reality.

#### 4.3.2 Hierarchal System

The system has one grid tree. The tree is constructed by connecting the GSS to 4 children where each child holds one backbone. Each backbone's root has 4 children where each child holds one Net. Each Net's root has two children where each child holds 5 sites.

#### 4.3.3 Hybrid System

The system has several grid trees acting as peers to each other. We use two hybrid systems in our simulation: 4 grid trees system (Hybrid-4T) and 16 grid trees system (Hybrid-16T). In the Hybrid-4T system the grid tree of the hierarchal system (see previous subsection) is broken to 4 grid trees where each backbone has one grid tree. In the Hybrid-16T (16 grid trees), each Net has one grid tree.

#### 4.4 Grid Jobs

A workstation submits a job to the grid via its grid entry (e.g. GSS) as a request that defines the job minimum requirements (JMR) for that job. Jobs are assumed to be executed until completion as batches on resources that meet their predefined requirements (i.e. operating systems, in our case). Gaussian distribution is usually used to simulate the required time to run a job on a server in [1, 9, 6, 10] with respect to the input job size and server's processing power. Therefore, we assume job sizes are correlated to the amount of work performed by each job where the input data size is expressed by Gaussian distribution with mean  $\mu = b * \text{cpus} * \text{wall time in seconds}$ , where  $b = 100$ , as in [9]. We also assume a job produces output data five times the original input job size.

#### 4.5 Computational Resources

We assume all resources (i.e. servers) are parallel machines that consist of a number of interconnected nodes with a number of CPUs within a node as in [1, 9]. Table 1 shows the servers used in the simulation experiments, which are originally based on real machines [1, 9, 6].

Server	Number of Nodes	CPUs per Node
1	184	16
2	305	4
3	144	8
4	1024	4
5	64	2
6	512	4
7	128	2

Table 1: Computational Resources

Those servers are duplicated in all the 160 sites in a range of 2 to 6 servers per site. The type of the servers and their operating systems (Windows, UNIX or LINUX) are picked up in random. We assume 0.1 local loads (i.e. not related to the grid) on all computational resources at all times, as the typical case in most studies like [6]. Note that the simulation model consists of 520 servers versus about 1880 workstations throughout all experiments— ratio 1:3.6.

#### 4.6 Workloads

Unfortunately it was difficult to find real traces for grid computing. However we were able to base our workloads on real traces for parallel machines and scientific applications [1, 6, 9, 12, 13]. Jobs in the workloads, that are relevant to this paper, use input average sizes of 1GB, 10GB and 100GB over the following number of jobs: 520, 1040, 1560, 3000 and 10,000 jobs. Refer to [1] for the complete set of workloads. We use Poisson distribution to generate input data sizes for submitted jobs to the grid where the Poisson mean is set to the desired average input size. In this way, jobs are generated with different sizes, but with the desired average input size, which is close to the typical case in reality.

#### 4.7 Performance Metrics

We use three performance metrics to compare systems: Total response time, average waiting time and average response time.

The total response time (TRT) is the time from submitting first job request until the completion of all jobs in a workload. For example, suppose that the first request was submitted to the grid at 5:00 PM and the last job of a workload was completed at 10:00 PM, the total response time will then be 5 hours. The purpose of this metric is to show the degree of parallelism in the grid, since we view the grid as a huge virtual parallel machine. The total response time (TRT) is calculated as follows:

$$TRT = (LJC - FRS)$$

where LJC (Last Job Completion Time) is the time that the output (i.e. at workstation) of the last completed job in the workload is received. FRS (First Request Submission), which is the time of transmitting the first request by a workstation.

The waiting time (WT) for a job is the time from submitting the job's request to the grid until the start of the actual job transfer to selected resources. For example, if a workstation submits a request to the grid at 5:00 PM, and gets a service offer from a resource at 6:00 PM, the waiting time for that job will be 1 hour. The purpose of this metric is to measure the scheduling time (i.e. the time it takes until a resource is allocated to that job). The average waiting time (AWT) is calculated as follows:

$$AWT = \frac{1}{N} \sum_{j=1}^N (SJTT - RT)_j$$

where  $N$  (Jobs count) is the number of jobs in a workload. SJTT (Start Job Transferring Time) is the time when a workstation receives a service offer from a resource and starts transferring the physical job. RT (Request Time) is the time when a workstation submits that job request to the grid.

The execution time (ET) is the time from submitting the actual job to the grid until the job's output is received at the

submitter's workstation. For example, a workstation receives a service offer from a resource at 5:00 PM. Suppose now that the workstation receives the output of the job at 6:00 PM, the execution time will then be 1 hour. Although all systems, in our model, function the same way when a request is mapped to a resource, but we still need this metric to measure the location where a job was executed at. The average execution time (AET) is calculated as follows:

$$AET = \frac{1}{N} \sum_{j=1}^N (JCT - SJTT)_j$$

where JCT (Job Completion Time) is the time when the job's output is received at the workstation.

#### 4.8 Simulation Experiments

We present, in this section, a sample of the experimental results in order to compare the performance over different scenarios. Refer to [1] for the complete set of results. Note that regardless of the configured system or experiment, the following assumptions still apply. (1) The computational power in the grid is maintained (i.e. 520 servers all the times). (2) A job is submitted by one workstation and executed by one server. Note that a workstation is called active if it has a pending request in the grid. Otherwise, it is called inactive. (3) A workstation that submits jobs according to a stochastic rate only operates at that rate while it is inactive. For example, a workstation submits jobs to the grid with rate 12 hours. Now, when that workstation becomes inactive, it waits according to that Poisson distribution with a mean of 12 hours before it submits another job. (4) All results are obtained by averaging 20 different runs. Note that the difference between the worse and the best case runs is in the range of 5% to 15%. Perhaps, this is because of having too many nodes in the model.

##### 4.8.1 First Experiment

In this experiment workstations submit jobs one after another until the entire workload is completed. This scenario is



possible when an organization, for instance, executes a number of jobs one after another automatically as a set. The workload in this experiment is already distributed among sites by the submission approach. For example, if site A has 3 workstations and site B has 6 workstations. Most likely, site B will submit twice the requests that are submitted from site A. The average waiting time (AWT) and the total response time (TRT) showed a substantial improvement against the P2P system regardless of the number of used grid trees, workload or scenario, as shown in Figures 4.1, 4.2, 4.4 and 4.7. Interestingly, the AWT starts declining when the hybrid system contains too many grid trees. Perhaps, this is because it gets closer and closer to the P2P system as a result of continue increasing the trees in the system. The average execution time (AET) is almost the same for small jobs (1 GB), but starts to differ when jobs-size increases (100 GB), as shown in Figures 4.3 and 4.5, which makes sense, since the model is built with high performance links.

#### 4.8.2 Second Experiment

In this experiment, workstations operate at different stochastic submission rate where each workstation selects, in random, one of the following rates: 10 minutes, 30 minutes, 1 hour, 5 hours, 1 day or 1 week. Furthermore, a workload in this experiment is not already distributed among sites, as in the case of the first experiment. Furthermore, in this scenario, sites also have different probabilities when generating a job. For example, if site A has 3 workstations and site B has 6 workstations. It is not necessary that site B is going to submit twice the requests that will be submitted from site A. On the other hand, it is quite possible that all requests will be submitted from site A. Observations of the first experiment are also supported by this experiment. Furthermore, both the butterfly and the load-balance algorithms showed a significant influence on the performance of the hierarchal system, as shown in Figure 4.6. In fact, the more jobs in a workload

the worse it gets, if the subject algorithms are disabled, hence, the more jobs, the more performed rescheduling.

#### 4.8.3 Third Experiment

In this experiment workstations submit jobs with the same stochastic rate (e.g. one hour rate for all workstations in the grid). A workload in this experiment is already distributed among sites, as in the case of the first experiment. We study the systems with three different rates: one hour, one day and one week.

Observations of the previous experiments are also supported by this experiment. In addition, the AWT tends to decline with a big slope in the hierarchal systems when jobs arrive into the grid with larger mean rate. However, it decreases slightly in the P2P system, as shown in Figure 4.7.

#### 4.8.4 Fourth Experiment

In this experiment, resources change according to one of the following stochastic changing-rates: One day, three days, one week, one month, three months and six months. Note that the changing-rate is also reselected, in random, along with advertised resources. For example, a server selects 6 months changing rate and reselects 3 months changing rate when it changes its advertised resources.

Now, when resources-change is a possibility during jobs scheduling, the AWT turns out the same comparing when resources are constant, as shown in Figure 4.8, which makes sense, since the fallback algorithm reschedules jobs while resources are already busy in executing other jobs anyway.

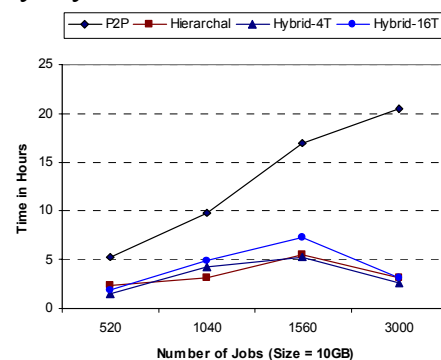


Figure 4.1: A Sample of Average Waiting Time in Experiment 1

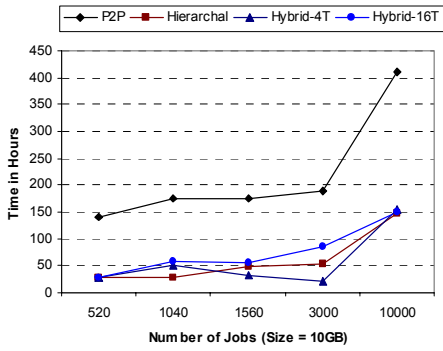


Figure 4.2: A Sample of Total Response Time in Experiment 1

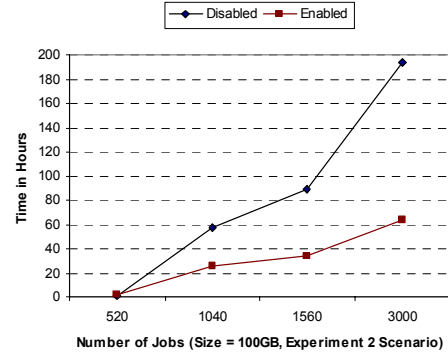


Figure 4.6: A Sample of algorithms influence on the Average Waiting Time in Experiment 2

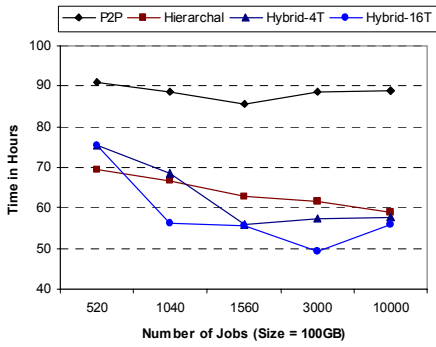


Figure 4.3: A Sample of Average Execution Time for large-size jobs in Experiment 1

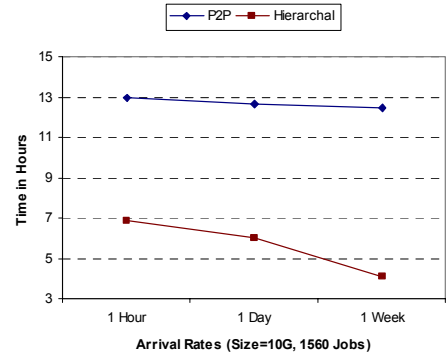


Figure 4.7: A Sample of Average Waiting Time in Experiment 3

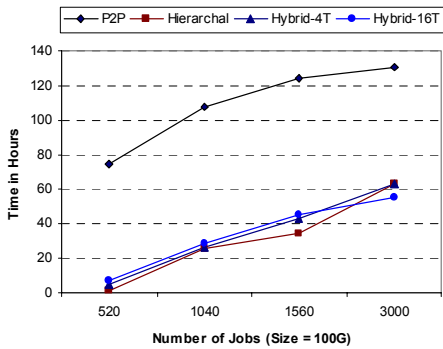


Figure 4.4: A Sample of Average Waiting Time in Experiment 2

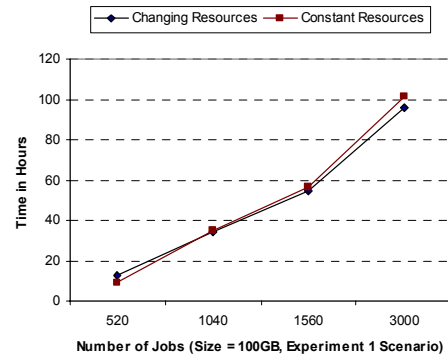


Figure 4.8: A Sample of Average Waiting Time in Experiment 4

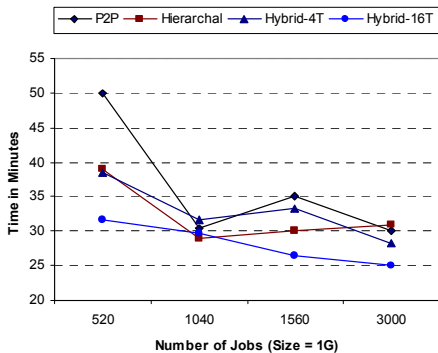


Figure 4.5: A Sample of Average Execution Time for small-size jobs in Experiment 2

## 5. Conclusions

Many studies jump over the resource discovery stage into the second scheduling stage by assuming that all jobs can execute anywhere in the grid, or simply assuming that resources will be discovered using the P2P approach. However, as we've shown those stages have to be dealt with in a sequence because of their dependency on each other. The hierarchal approach has not only shown substantial improvement over the P2P system, but also the ability to

be combined with it in one hybrid system. Both the average waiting time (AWT) and the total response time (TRT) metrics showed a large improvement of the hierarchical approach against the P2P system regardless of the number of used grid trees, workload or scenario. The average execution time (AET) metric also showed a significant improvement over the P2P system for large-size jobs, but almost the same for small-size jobs, which makes sense, since the model is built with high performance links. The three rescheduling algorithms showed a big contribution in the overall performance of the system. The fallback algorithm saved not only some jobs of not being able to execute because of resources change, but also maintained the same system performance when resources are constant. Both the butterfly and load-balance algorithms saved the system of performing poorly when increasing the number of jobs in workloads.

### References:

- [1] K. AlZoubi, *Hierarchical Scheduling in Grid systems*. Master Thesis, Carleton University, Ottawa, Canada, 2006 (expected).
- [2] V. Berstis, “Fundamentals of Grid Computing”, <http://www.redbooks.ibm.com/redpapers/pdfs/redp3613.pdf>
- [3] S. Dandamudi, *Hierarchical Scheduling in Parallel and Cluster systems*. Kluwer Academic Publishers, 2003.
- [4] I. Foster and C. Kesselman, *The Grid: Blueprint for a New Computing Infrastructure*, Morgan Kaufmann, 2004.
- [5] I. Foster, C. Kesselman, and S. Tuecke, “The Anatomy of the Grid: Enabling Scalable Virtual Organizations”, <http://www.globus.org/research/papers/anatomy.pdf>.
- [6] S. Hotovy, “Analysis of the Early Workload on the Cornell Theory”, ACM SIGMETRICS, 1996, pp. 272 – 273.
- [7] J. Nabrzyski, J. Schopf and J. Weglarz, *Grid Resource Management: State of the Art and Future Trends*. Kluwer Academic Publishers, 2004.
- [8] A. Odlyzko, “Internet Traffic Growth: Sources and Implications”, <http://www.dtc.umn.edu/~odlyzko/doc/itcom.internet.growth.pdf>
- [9] H. Shan, W. Smith, L. Oliner and R. Biswas, “Job Scheduling in a Heterogeneous Grid Environment”, <http://www-library.lbl.gov/docs/LBNL/549/06/PDF/LBNL-54906.pdf>
- [10] A. Takefusa, “Bricks: A Performance Evaluation System for Scheduling Algorithms on the Grids”, JWAITS 2001.
- [11] G. Wainer. *DEVS CD++ tools*, <http://www.sce.carleton.ca/faculty/wainer.html>.
- [12] Parallel Workloads Archive. <http://www.cs.huji.ac.il/labs/parallel/workload/> .
- [13] “Long Term Technology Review of the Science & Engineering Base”, [http://www.rcuk.ac.uk/ltr/finalversion/LTTR\\_ContentPage.htm](http://www.rcuk.ac.uk/ltr/finalversion/LTTR_ContentPage.htm), 2000.