

# A Metascheduler For The Grid\*

Sathish S. Vadhiyar  
Computer Science Department  
University of Tennessee  
vss@cs.utk.edu

Jack J. Dongarra  
Computer Science Department  
University of Tennessee  
dongarra@cs.utk.edu

## Abstract

*With the advent of Grid Computing, scheduling strategies for distributed heterogeneous systems have either become irrelevant or have to be extended significantly to support Grid dynamics. In this paper, we describe a metascheduling architecture for a Grid system that takes into account both the application and system level considerations. Results are presented to demonstrate the usefulness of the metascheduler.*

## 1. Introduction

There have been number of efforts in devising and/or implementing scheduling strategies for heterogeneous distributed computing systems since the advent of Network of Workstations (NOWs) [5], [14], [21], [8], [17], [18], [2]. The Grid [6] is an abstraction of distributed heterogeneous systems and investigating the relevance of scheduling strategies for heterogeneous systems to Grid environment is a worthwhile effort. The work by Khaled Al-Saqabi et. al [14] considers a 2D array of processors and time slices and assigns the Virtual Processes (VPs) of the jobs to the array. Scheduling based on time slices will lead to huge overhead for the scheduling system when the scheduling strategies have to be invoked frequently in response to frequent Grid dynamics. The Load Sharing Facility [21] lays emphasis on distributing the jobs among the available machines based on the workload on the machines. The assumption that load sharing leads to good response times is not valid in a Grid scenario where the network heterogeneity can significantly affect the execution time of the application.

MARS [8] and more recently AppLeS [2] provide good approaches for application level scheduling in meta computing environments. AppLeS is more suitable for Grid environment with its sophisticated NWS [19] mechanism for collecting system information. However, both MARS and

AppLeS do not have powerful resource managers that can negotiate with applications to balance the interests of different applications. The absence of these negotiating mechanisms in a Grid can lead to various problems like the bushel of AppLeS problem [2].

In this paper, we describe a metascheduling architecture that we have been building in the context of the GrADS project [1]. The metascheduler receives candidate schedules of different application level schedulers and implements scheduling policies for balancing the interests of different applications. The goals of the metascheduler include:

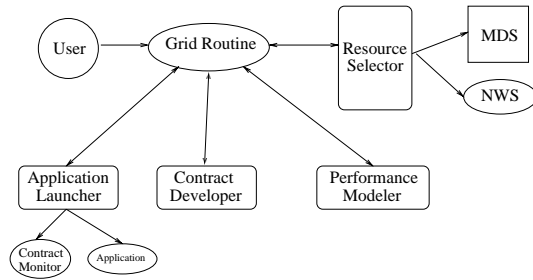
1. Verifying that the applications made their scheduling decisions based on conditions of the system when competing applications are executing.
2. Accommodating short running jobs by temporarily stopping long running and resource consuming jobs.
3. Facilitating new applications to execute faster by stopping certain competing applications.
4. Minimizing the impact that new applications can create on already running applications.
5. Migrating running applications to new machines in response to system load changes to improve the performance or to prevent performance degradation.

In Section 2, we give a brief overview of the existing GrADS project that utilizes application level scheduling. In Section 3, we describe the metascheduler that we have been building for GrADS environment. We explain in detail the different components of the metascheduler and the mechanisms in the components to achieve the goals mentioned above. In Section 4, we present experiments and results to validate the usefulness of the metascheduler. In Section 5, we compare our metascheduler with related efforts. In Section 6, we present some conclusions. In Section 7, we mention some of the future plans for our metascheduler.

---

\*This work is supported in part by the National Science Foundation contract GRANT #E81-9975020, SC R3605-29200099, R01-1030-09.

**Figure 1. GrADS Architecture for Numerical Libraries**

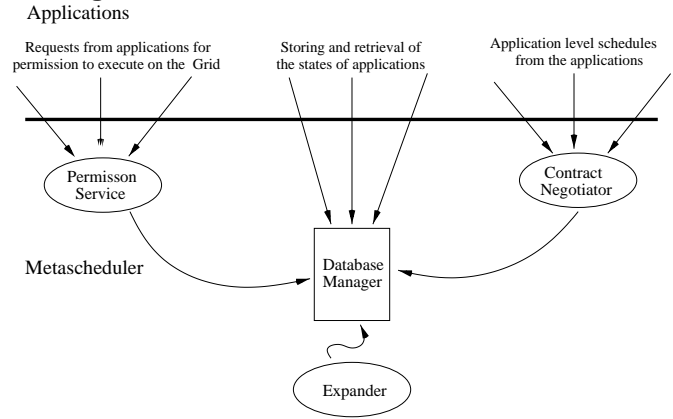


## 2. The GrADS system

GrADS [1] is an ongoing research involving number of institutions and its goal is to simplify distributed heterogeneous computing in the same way that the World Wide Web simplified information sharing over the Internet. The University of Tennessee investigates issues regarding integration of numerical libraries in the GrADS system. In our previous work [11], we demonstrated the ease in which numerical libraries like ScaLAPACK can be integrated into the Grid system and the ease in which the libraries can be used over the Grid. We also showed some results to prove the usefulness of a Grid in solving large numerical problems. The architecture that was used in the work is illustrated by Figure 1.

As a first step, the user invokes a Grid routine with the problem he wants to solve along with the problem parameters. The Grid routine invokes a component called Resource Selector. The Resource Selector accesses the Globus MetaDirectory Service(MDS) to get a list of machines that are alive and then contacts the Network Weather Service(NWS) to get system information for the machines. The Grid routine then invokes a component called Performance Modeler with problem parameters, machines and machine information. The Performance Modeler through an execution model built specifically for the application, determines the final list of machines for application execution. By employing the application specific execution model, GrADS follows the AppLeS approach to scheduling. The problem parameters and the final list of machines are passed as a contract to a component called Contract Developer. The Contract Developer is primitive in that it approves all the contracts that are passed to it. The Grid routine then passes the problem, its parameters and the final list of machines to Application Launcher. The Application Launcher spawns the job on the given machines using Globus job management mechanism and also spawns a component called Contract Monitor. The Contract Monitor through an Autopilot mechanism [13] monitors the times

**Figure 2. Metascheduler and interactions**



taken for different parts of applications and displays the actual and predicted times. Eventually the Contract Monitor will be used for sending information about contract violations to a rescheduler which can in turn take corrective measures on the application execution.

## 3. Metascheduling in the GrADS architecture

The architecture shown in Figure 1 implements application level scheduling through the use of execution model built specifically for the application. The execution model does not take into account the existence of other applications in the system. There are a number of potential problems with the application level scheduling implemented by the architecture. First, when two applications are submitted to the Grid at the same time, scheduling decisions will be made for each application assuming the absence of the other application. Second, in the above architecture, if, through the performance model, a new job submitted to the Grid system detects that the Grid resources are not sufficient for it to execute, it cannot make further progress. Similarly, a long running job that was submitted to the system can severely impact the performance of new jobs that enter the system. The root cause of the above and other problems is the absence of a metascheduler that obtains the candidate schedules from different applications and try to balance the needs of different applications. The metascheduler is implemented by the addition of four new components, namely, database manager, permission service, contract negotiator and expander, to the architecture shown in Figure 1. The interactions between these different metascheduler components and the interactions between the applications and the metascheduler are illustrated by Figure 2.

The following subsections describe each of the metascheduler components.

### 3.1. Database Manager

The database manager maintains a record for each application submitted to the Grid system. The record contains the state of the application, the resource information of the Grid resources when the application entered the system, the final list of machines on which the application executes, the predicted time for the application etc. This information is queried by the other components of the metascheduler to make scheduling decisions.

### 3.2. Permission Service

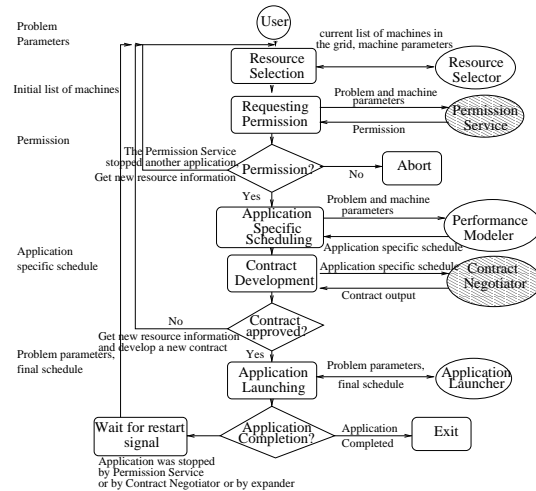
Permission Service is a daemon that receives requests from the applications to grant them permission to proceed with the usage of the Grid system. The Permission service checks if the Grid resources have adequate free memory to execute the application. If the free memory of the Grid resources is less than the free memory required by the application, then executing the application on the resources will lead to large execution time of the application due to frequent access to local disks. Hence in this case, the Permission Service either denies permission for the request or tries to accommodate the application by stopping an already executing resource and time consuming application. The functions of the Permission Service are illustrated by the pseudo code in Appendix A.

### 3.3. Contract Negotiator

The Contract Negotiator component of the metascheduler is a daemon that receives application level schedules from the applications. An application level schedule of an application is the final list of machines that the application obtains from the Performance Modeler through the employment of the application specific execution model. These are the list of machines on which the application can potentially execute. The application passes the problem parameters and the application level schedule in the form of a contract to the Contract Developer. The Contract Developer, instead of approving the contracts of the applications under all conditions, contacts the Contract Negotiator for obtaining approval of the application contract. The Contract Negotiator acts as a queue manager controlling different applications of the Grid system. The Contract Negotiator either approves the contract in which case the application can proceed to the application launching phase, or rejects the contract in which case the application restarts from the resource selection phase. The Contract Negotiator rejects the contract under the following conditions:

1. When the application has got its resource information from NWS before an executing application started executing.

Figure 3. Life cycle of an application in the Grid



2. If the performance of the new application can be improved significantly in the absence of an executing application. An executing application may have large remaining execution time and may be consuming large number of resources, the availability of which can significantly enhance the performance of the new application. In this case, the contract negotiator either waits for the executing application to complete or proactively stops the executing application to accommodate the new application.
3. If the already executing applications can be severely impacted by the new application.

1 and 2 have already been implemented while 3 is a work in progress. The working of the Contract Negotiator is illustrated by the pseudo code in Appendix B.

### 3.4. Expander

Expander is a daemon that tries to improve the performance of the already executing applications. It queries the database manager at regular intervals for completed applications. When an application completes, the expander determines if performance benefits can be obtained for an already executing application by expanding the application to utilize the resources freed by the completed application. If the expander detects such an executing application, it stops the application and continues the application on the new set of resources. The pseudo code of the expander is given in Appendix C.

The life cycle of an application and its interactions with the metascheduler is shown in Figure 3.

**Table 1. Machine specifications**

Machine name	Processor type	Speed (MHz)	Memory (MByte)	Network
torc	Pentium III	550	512	100 Mb switched Ethernet
msc	Pentium III	933	512	100 Mb switched Ethernet
cypher	Pentium III	500	512	1 Gbit switched Ethernet
opus	Pentium II	450	256	1.28 Gbit/sec full duplex myrinet

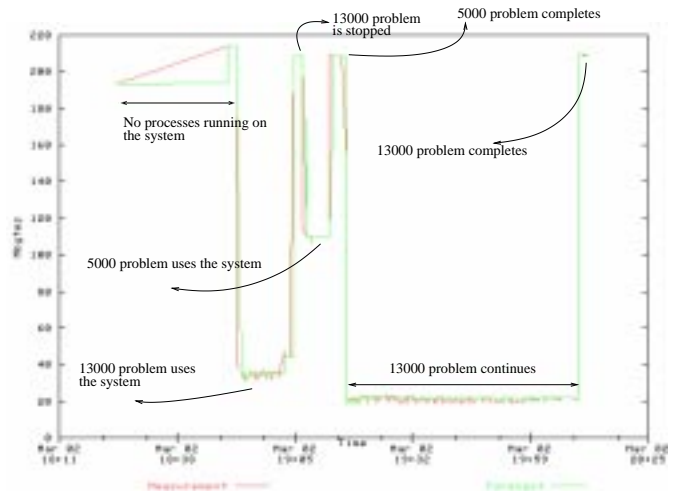
## 4. Experiments and Results

ScaLAPACK LU and QR factorization codes were instrumented such that the time taken for each iteration corresponding to a block of the matrix is measured and monitored. Mechanisms have been implemented in the ScaLAPACK code that will enable the ScaLAPACK application to be stopped and restarted on possibly different number of processors. We use the Internet Backplane Protocol (IBP) [12] for storage of the checkpoint states. IBP depots, where storage can be allocated, are started on the processors of the Grid System.

The GrADS experimental testbed consists of about 40 machines that reside in institutions across the country including University of Tennessee, University of Illinois, University of California at San Diego, Rice University etc. For the easy demonstration of our experimental results, our experimental testbed consists of a cluster in UIUC called *opus* consisting of 8 machines, a cluster in University of Tennessee called *torc* consisting of 8 machines, another cluster in University of Tennessee called *msc* consisting of 8 machines and another cluster in University of Tennessee called *cypher* consisting of 16 machines. The *opus* cluster is connected to the other 3 clusters clusters by Internet. *torcs*, *mscs* and *cyphers* are connected to each other 100 Mb Ethernet links. Table 1 gives the specification of the machines.

The total execution times reported in the following subsections include the time for Grid overhead and not just the time taken by the actual application. The time for the Grid overhead is reported in our previous work [11].

**Figure 4. Free memory available on a opus machine during the execution of  $app_1$  and  $app_2$**

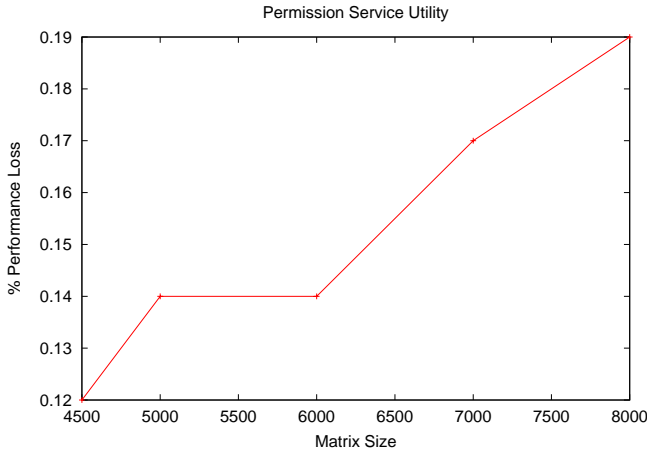


### 4.1. Experiment 1

In this experiment, we demonstrate the functionality of the Permission Service. For the experiments in this section, ScaLAPACK LU factorization code was used. A large application,  $app_1$ , was introduced into the system consisting of 4 *opus* machines, 1 *torc* machine and 2 *cypher* machines. Ten minutes after  $app_1$  started, a relatively small application,  $app_2$ , that intended to use only the 4 *opus* machines was introduced into the system.  $app_2$  was chosen such that its memory requirements were greater than the memory available in the *opus* system when  $app_1$  was executing. In the following experiment, a linear algebra problem with matrix size 13000 was chosen for  $app_1$ . The Permission Service evaluated the performance benefits of stopping  $app_1$ , accommodating  $app_2$ , and restarting  $app_1$  after the completion of  $app_2$ . The functionality of the Permission Service, when the matrix size of the linear algebra problem,  $app_2$ , is 5000, is illustrated on a single *opus* machine in Figure 4. The graph was generated in the NWS web site.

In Figure 5, we observe the percentage performance loss incurred by  $app_1$  due to the accommodation of  $app_2$  in the system. The x-axis represents different matrix sizes for  $app_2$  and the y-axis represents the percentage performance loss incurred by  $app_1$ . Two points can be observed from Figure 5. First, for less than 20% of performance loss for  $app_1$ , the system was able to accommodate  $app_2$ . Without the Permission Service mechanism,  $app_2$  would not have been able to use the system. Second, the performance loss increases with the increase in problem size of  $app_2$ . When the problem size of  $app_2$  is comparable with the problem size of  $app_1$ , the Permission Service determines that perfor-

**Figure 5. Performance loss for  $app_1$**



mance benefits cannot be achieved for the system by accommodating  $app_2$ . In order to prevent continuous preemption of  $app_1$  by small applications, the scheduling strategy is implemented such that  $app_1$  is ensured to make at least 20% of progress between preemptions.

## 4.2. Experiment 2

In this experiment, we demonstrate the utility of the contract negotiator in accommodating a new application by stopping an already running application, if significant performance benefits can be obtained for the new application. The stopped application is restarted after the new application completes its execution. For this experiment, ScaLAPACK LU factorization code was used on only *cypher* machines. In this experiment, an application,  $app_1$  is executed on  $N$  processors. 3 minutes after  $app_1$  started its execution, an application,  $app_2$  is introduced in the Grid system.  $app_2$  is intended to use  $(N+1)$  processors. Since  $N$  of the processors were occupied by  $app_1$ , only a single processor is available for  $app_2$ . The Contract Negotiator analyzes the performance benefits that can be obtained by stopping  $app_1$  and making  $(N+1)$  processors available for  $app_2$ . In the experiments, matrix size 7500 was used for  $app_2$ . The total execution time of a 7500 matrix size ScaLAPACK problem when executed on a single processor is 818.11 seconds.

We define

1. Execution time of  $app_1$  without rescheduling,  $exec1_{without\_re}$
2. Execution time of  $app_1$  with rescheduling,  $exec1_{with\_re}$
3. Execution time of  $app_2$  without rescheduling,  $exec2_{without\_re}$

**Table 2. Utility of Contract Negotiator**

Matrix Size of $app_1$	Processors $N$	Number of Processors used by $app_2$	util_val
15000	4	5	2.13
17000	5	6	5.11
18500	6	7	2.27
20000	7	8	2.04
21000	8	9	2.05
22500	9	9	2.36
24000	10	9	1.72

4. Execution time of  $app_2$  with rescheduling,  $exec2_{with\_re}$

5. Performance loss for  $app_1$ , perf\_loss

$$perf\_loss = \frac{exec1_{with\_re} - exec1_{without\_re}}{exec1_{without\_re}}$$

6. Performance gain for  $app_2$ , perf\_gain

$$perf\_gain = \frac{exec2_{without\_re} - exec2_{with\_re}}{exec2_{without\_re}}$$

7. Utility value, util\_val

$$util\_val = \frac{perf\_gain}{perf\_loss}$$

$util\_val > 1$  indicates that the rescheduling strategy is useful for the entire system.  $util\_val < 1$  indicates that the rescheduling strategy can cause an overall loss in performance for the entire system. Greater the value of util\_val, more the usefulness of the rescheduling strategy.

Table 2 shows the matrix sizes of  $app_1$ , the number of processors  $N$ , the number of processors eventually used by  $app_2$  and the util\_val. Note that the eventual number of processors used by  $app_2$  depends on system conditions and execution time model and is not always the  $(N+1)$  processors available to  $app_2$ .

We observe from Table 2, that the values of util\_val are consistently high for the above experiments. This showed that the scheduling strategy of compromising long running jobs for short running jobs is beneficial to the entire system. The value of util\_val depends on a number of factors including the times for the long and short jobs and the times for checkpointing the states of the long job. As in the Permission Service, mechanisms have been implemented to avoid continuous preemptions of  $app_1$ .

### 4.3. Experiment 3

In this set of experiments, we illustrate the utility of Expander. For the experiments in this section, ScaLAPACK QR factorization code was used. An application,  $app_1$ , was introduced into the system such that it consumed most of the memory of 8  $msc$  machines. During the execution of  $app_1$ , an  $app_2$ , that intended to use 11 machines, 3  $torcs$  and 8  $mscs$  was introduced into the system. Since the 8  $msc$  machines were occupied by  $app_1$ ,  $app_2$  was able to utilize only the 3  $torc$  machines. When  $app_1$  completed, the 8  $msc$  machines were freed and  $app_2$  was able to utilize the extra resources to reduce its remaining execution time. The Expander evaluated the performance benefits of allowing  $app_2$  to utilize the extra 8 processors.

ScaLAPACK problems of sizes 20000 and 21000, depending on the available memory on  $mscs$  when the experiments were run, were used for  $app_1$ . ScaLAPACK problem of size 11000 was used for  $app_2$ .

We define

1. Total execution time of  $app_2$  on 3  $torcs$  without rescheduling,  $exec_{without\_re}$
2. Total execution time of  $app_2$  with rescheduling,  $exec_{with\_re}$
3. Percentage rescheduling gain for  $app_2$ ,  $percentage\_gain$

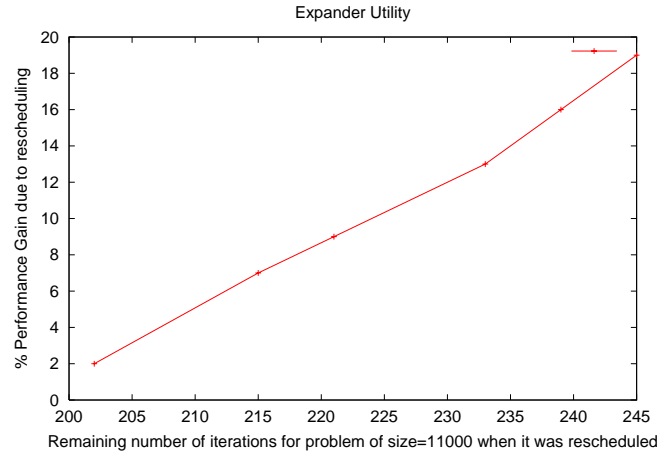
$$percentage\_gain = \frac{exec_{without\_re} - exec_{with\_re}}{exec_{without\_re}}$$

$app_2$  was introduced at various points of time after the starting of  $app_1$ . Hence additional resources will be available for  $app_2$  at various points of time into its execution. The total number of iterations needed by the ScaLAPACK problem of size 11000 was 275. Figure 6 illustrates the utility of rescheduling as a function of the remaining number of iterations left for  $app_2$  when  $app_2$  was rescheduled. We observe that the percentage rescheduling gain for  $app_2$  increases when the remaining execution time left for  $app_2$  at the time of rescheduling increases. The rescheduling gain depends on a number of parameters like the time taken for redistribution of data and the number of additional resources available etc. These parameters depend on the specific application for which rescheduling is done. Work is in progress to build interfaces for the application library writer by which the system can determine the rescheduling parameters for the applications.

### 5. Related Work

There are number of ongoing research efforts in Grid Computing [7], [9], [10], [4], [3], [16]. There are

Figure 6. Rescheduling gain for  $app_2$



number of similarities in the scheduling systems of Condor [10] and our metascheduler. Condor also supports preemption of executing jobs to either accommodate other jobs [20] or to transfer the control of the resources to the resource owners. Moreover, on a broad level, the functionality of our Contract Negotiator in our metascheduler is similar to the functionality of the *negotiator* in the Condor system in that these negotiator components negotiate between the applications and the resources. But the differences between the Condor scheduling system and our metascheduler lie in both the overall objectives and the capabilities of the systems. While the main motivation of the scheduling decisions in Condor is to utilize idle resources, the motivation for our metascheduler is to improve the performance of the individual applications. In Condor, jobs are preempted from utilizing the resources when the resources are reclaimed by the owners, whereas in the metascheduler, jobs are preempted from resources when the availability of the resources can significantly improve the performance of other jobs. Moreover, at present, Condor does not support checkpointing of parallel jobs. The objectives of Nimrod-G's [3] scheduling policies are similar to those of our metascheduler where different users' requirements are balanced. Nimrod-G uses grid economies to implement its scheduling policies while our metascheduler uses predicted application times for our scheduling policies. Though the Ninf [16] team had evaluated their scheduler when multiple clients run their jobs, no substantial mechanism has been implemented to guarantee performance for each client. The metascheduler adheres to the definition of the Super Scheduler described in the Global Grid Forum working document [15]. In addition, the metascheduler acts as a Super Scheduler that negotiates the decisions made by the different Super Schedulers described in the working document.

## 6. Conclusion

In this paper we have explained the implementation of a metascheduler for the Grid that takes into account both the application level and system level considerations. We have explained in detail, the different components of our metascheduler, viz., the Permission Service, Contract Negotiator and the Expander. These components provide valuable scheduling services that play important roles in providing scalability of the Grid system. We have demonstrated the utility of these scheduling decisions with encouraging results.

## 7. Future Work

Our immediate plans are to demonstrate the usefulness of our metascheduler in a large Grid system when large number of applications execute on the system. Capabilities like evaluating the impact of new applications on existing applications and migration under performance degradation will be added to the metascheduler. After the complete implementation of the metascheduler, the issue of reproducibility of numerical results in the Grid will be investigated.

## References

- [1] F. Berman, A. Chien, K. Cooper, J. Dongarra, I. Foster, D. Gannon, L. Johnsson, K. Kennedy, C. Kesselman, J. Mellor-Crummey, D. Reed, L. Torczon, and R. Wolski. The GrADS Project: Software Support for High-Level Grid Application Development. *International Journal of High Performance Applications and Supercomputing*, 15(4):327–344, Winter 2001.
- [2] F. Berman and R. Wolski. The AppLeS Project: A Status Report. *Proceedings of the 8th NEC Research Symposium*, May 1997.
- [3] R. Buyya, D. Abramson, and J. Giddy. Nimrod-G Resource Broker for Service-Oriented Grid Computing. *IEEE Distributed Systems Online*, 2(7), November 2001.
- [4] H. Casanova and J. Dongarra. NetSolve: A Network Server for Solving Computational Science Problems. *The International Journal of Supercomputer Applications and High Performance Computing*, 11(3):212–223, Fall 1997.
- [5] T. Casavant and J. Kuhl. A Taxonomy of Scheduling in General-Purpose Distributed Computing Systems. *IEEE Transactions on Software Engineering*, SE-14(2):141–154, February 1988.
- [6] I. Foster and C. K. eds. *The Grid: Blueprint for a New Computing Infrastructure*. Morgan Kaufmann, ISBN 1-55860-475-8, 1999.
- [7] I. Foster and C. Kesselman. Globus: A Metacomputing Infrastructure Toolkit. *Intl J. Supercomputer Applications*, 11(2):115–128, 1997.
- [8] J. Gehring and A. Reinefeld. MARS - A Framework for Minimizing the Job Execution Time in a Metacomputing Environment. *Future Generation Computer Systems*, 12(1):87–99, 1996.
- [9] A. Grimshaw, W. Wulf, J. French, A. Weaver, and J. P. Reynolds. Legion: The Next Logical Step Toward a Nationwide Virtual Computer. Technical Report CS-94-21, Department of Computer Science, University of Virginia, 1994.
- [10] M. Litzkow, M. Livney, and M. Mutka. Condor - a Hunter for Idle Workstations. *Proc. 8th Intl. Conf. on Distributed Computing Systems*, pages 104–111, 1988.
- [11] A. Petitet, S. Blackford, J. Dongarra, B. Ellis, G. Fagg, K. Roche, and S. Vadhiyar. Numerical Libraries and the Grid: The Grads Experiments with Scalapack. *Journal of High Performance Applications and Supercomputing*, 15(4):359–374, Winter 2001.
- [12] J. S. Plank, M. Beck, W. R. Elwasif, T. Moore, M. Swany, and R. Wolski. The Internet Backplane Protocol: Storage in the Network. *NetStore99: The Network Storage Symposium*, 1999.
- [13] R. Ribler, J. Vetter, H. Simitci, and D. Reed. Autopilot: Adaptive Control of Distributed Applications. *Proceedings of the 7th IEEE Symposium on High-Performance Distributed Computing*, July 1998.
- [14] K. Saqabi, S. Otto, and J. Walpole. Gang Scheduling in Heterogeneous Distributed Systems. Technical report, OGI, 1994.
- [15] J. Schopf. Super Scheduler Steps/Framework. <http://www-unix.mcs.anl.gov/~schopf/ggf-sched/>, July 2001.
- [16] S. Sekiguchi, M. Sato, H. Nakada, and U. Nagashima. Ninf: Network based Information Library for Globally High Performance Computing. *Parallel Object-Oriented Methods and Applications (POOMA)*, February 1996.
- [17] C. Waldspurger and W. Weihl. Lottery Scheduling: Flexible Proportional-Share Resource Management. *First Symposium on Operating Systems Design and Implementation (OSDI)*, pages 1–11, 1995.
- [18] J. Weissman. The Interference Paradigm for Network Job Scheduling. *Proceedings of the Heterogeneous Computing Workshop*, pages 38–45, April 1996.
- [19] R. Wolski, N. Spring, and J. Hayes. The Network Weather Service: A Distributed Resource Performance Forecasting Service for Metacomputing. *Journal of Future Generation Computing Systems*, 15(5-6):757–768, October 1999.
- [20] D. Wright. Cheap Cycles from the Desktop to the Dedicated Cluster: Combining Opportunistic and Dedicated Scheduling with Condor. *Proceedings of the Linux Clusters: The HPC Revolution conference, Champaign - Urbana, IL*, June 2001.
- [21] S. Zhou, X. Zheng, J. Wang, and P. Delisle. Utopia: a Load Sharing Facility for Large, Heterogeneous Distributed Computer Systems. *Software - Practice and Experience*, 23(12):1305–1336, December 1993.

## A Appendix A - Permission Service

```
Input: problemParameters,
       resourceRequirements,
       currentResourceInformation
resourceCapacity =
  getCapacity
  (currentResourceInformation)
if (resourceCapacity >
    resourceRequirements)
  send (PERMISSION)
else
  largeApplications =
    get list of applications
    whose memory requirements
    are greater than the
    memory requirement for
    this application
  if (largeApplications is empty)
    send (NO_PERMISSION)
  else
    shortRemainingTimeAppl =
      resource consuming applications
      that are going to end in 3
      minutes
    if (shortRemainingTimeApp) then
      waitForCompletion
      (shortRemainingTimeApp)
      send(PERMISSION)
    else
      ratio = remaining execution
              time of
              largeApplication /
              predicted execution
              time of new application
              in the absence of
              largeApplication
      maxRatio = maximum value of
                  ratio for all
                  largeApplications
      if (maxRatio >20 and
          largeApplication is
          checkpointEnabled)
        stop largeApplication
        send permission to new
        application
        wait for new application to
        complete
        resume largeApplication
```

## B Appendix B - Contract Negotiator

```
Input: problemParameters,
       finalListOfMachines,
       predictedTime
rsTime = time of resource selection
         for this application
executingList =
  getListOfExecutingApplications()
if (rsTime < minimum starting time
    of applications in
    executingList)
  send (CONTRACT_NOT_OK)
else
  for each application i in
  executingList
    timeAbs = predicted time of new
              application in the
              absence of i
    timePre = predicted time of new
              application in the
              presence of i
    if (timePre/timeAbs > 1.5)
      bigApplication = i
      break out of for each loop
    remainingExecAbs =
      remaining execution time of
      application i in the absence of
      new application
    remainingExecPre =
      remaining execution time of
      application i in the presence
      of new application
    impactTime = remainingExecPre -
                 remainingExecAbs
    if (application i is
        checkpointable and
        2*timeAbs < 0.5*
        min(remainingExecTime+timeAbs,
            impactTime + timePre))
      stop application i;
      send (CONTRACT_NOT_OK) to new
      application
      wait for the new application to
      complete
      resume application i
    else if ((impactTime + timePre)
             > 1.2*(remainingExecTime+timeAbs))
      Continue application i
      Wait for application i to complete
      Send (CONTRACT_NOT_OK) to new
      application
    else
      Send (CONTRACT_OK) to new
      application
```



## C Appendix C - Expander

```
recentlyCompletedAppl =
  applications that completed a minute
  ago
if (recentlyCompletedApp)
  executingList =
    get list of executing applications
  subsetList =
    subset of executingList that
    are well behaved, i.e., whose
    actual performance is close
    to predicted
  for each application i in subset
    reschedulingGain =
      (remaining execution time of
      i without rescheduling -
      (remaining execution time of
      i with rescheduling +
      rescheduling time)
      ) /
      remaining execution time of i
      without rescheduling
  maxReschedulingGain =
    maximum of rescheduling gains
  maxApplication =
    application that has maximum
    rescheduling gain
  if (maxReschedulingGain > 0.5)
    stop maxApplication
    get new candidate schedule for
    maxApplication
    continue maxApplication on new
    set of resources
```