

Abstract

!!!TODO Skrive om abstract til slutt!!! This project looks at ways to minimize scheduling and administration overhead for very short independent tasks on heterogeneous grids. Dagoc, a commercial application for the Oil/-Gas industry, is used as the test case. Each task is between 1 and 30 seconds long and are currently calculated serially on one processor. Splitting such small tasks into individual jobs might not be suitable for grid systems. Ways to collect multiple tasks into "multi task jobs" are suggested to see if gives any speedup. Each task uses a set of input data files which resides on a NFS server which needs to be taken into the equation when distributing all the tasks.

TODO: What was achieved? Conclusion? Benchmarks gave this and that result...

Contents

1	Introduction	4
2	Related Work	5
3	Grids	6
3.1	Grid features	7
3.1.1	Condor 6.8.6	7
3.1.2	Sun Grid Engine 6 (SGE)	8
4	Implementation Ideas	9
4.1	A golden rule	9
4.2	Scheduler tuning	10
4.3	Submitting multiple jobs with one submit file	10
4.4	Bash script with multiple executions	11
4.5	Condor specific tools	11
4.5.1	DAGMan	11
4.5.2	MW: Master-Worker	11
4.6	Handling the input/output files for Dagoc	12
4.6.1	Condor specific file handling	12
4.6.2	Result gathering with DAGMan	13
4.6.3	Result gathering with SGE	13
4.6.4	Problems with NFS file access	13
4.7	Other areas of interest	14
5	Model	14
5.1	Test parameter	15
6	Benchmarks and Results	16
6.1	Serial execution	17
6.2	Results for Condor	17
6.3	Results for Sun Grid Engine (without result file transfer)	20
6.4	Results Sun Grid Engine (with result file transfer)	21
6.5	Results when altering the Job-Task ratio	23
7	Conclusion	25
7.1	Future work	27
A	Setting up the grid	29
A.1	Sun Grid Engine 6.1	29
A.1.1	Possible node configurations	29
A.1.2	Access to files	30
A.1.3	Problems	30
A.1.4	Windows restrictions	30

A.2	Condor-6.8.6	30
A.2.1	Possible node configurations	31
A.2.2	Access to files	31
A.2.3	Heterogeneous nodes	31
A.2.4	Problems	32
A.2.5	Windows restrictions	32
B	Using the grid	32
B.1	Sun Grid Engine 6.1	32
B.1.1	Submitting Jobs	32
B.2	Condor-6.8.6	33
B.2.1	Preparing jobs for execution	33
B.2.2	Submitting Jobs	33
B.2.3	Different settings	34
C	Scripts	34

1 Introduction

In recent years, multicore and parallel processing has become the norm for computer architectures. To utilize these new architectures to the full, software has to be (re)designed to do so. A serial code cannot utilize more than one core on a multicore processor. However, multiple independent serial applications will utilize multiple cores, if executed at the same time. Serial code can also be rewritten to make use of multiple cores using ,e.g., OpenMP or plain threads. However, this will not only require a lot of work, but knowledge about parallel programming paradigms is also required by the software developers. Thus, even more time and money must be spent gaining this knowledge. Eventually, gathering this knowledge will become a necessity, as multicore is the future.

In this project, we consider a commercial application, Dagoc, which is a tool for the Oil/Gas industry. It is developed by a small, Norwegian software company called Yggdrasil AS. The part of the application we are looking is a typical Parameter Sweep Application (PSA) [10]. We consider the computation of different sized collections of small, equal-sized and independent tasks. The collections are currently designed to be executed serially on a single workstation. Each task in a collection computes a fixed amount of source data from a list of different input parameters. The final result of a collection is to be gathered at the end. At first look, our application should be perfect for use in a Grid environment.

However, the individual tasks are quite small, between 2 and 20 seconds. An average task needs about 0.5MB of input data from multiple small files. The executable is about 22MB and needs about 1.5MB of custom libraries. If the program is to be executed in database mode, it will need access to an 1-2MB sqlite3 database file as well. For such short tasks, the scheduling and file distribution may become a dominant factor.

The main focus of this project will be looking at how to apply Grid technology to support the jobs in our application. Using Grid technology, multiple cores can be utilized, even multiple computers with multiple cores in parallel. Methods how to minimize scheduling overhead using built in Grid middleware functions are also considered. These methods do not include altering the source code in the application to a great extent.

A typical office environment is considered. Here, the collection of workstations will in most cases be quite heterogeneous. Some machines are upgraded while others are kept as is. Some nodes might use Windows or Mac OS X while others use different flavors of Linux. Binary compatibility between these platforms will induce problems. Sometimes between different flavors of Linux as well. An executable compiled for one flavor of Linux might return an error on another, because of missing or outdated libraries or other causes. A Grid must be able to handle this. If jobs are distributed carelessly the success of each job execution could end up being quite ran-

dom. Without any constraints, the Grid engine will send jobs to arbitrary free nodes. The user will not know if their job will be executed successfully or if it will return an error. A common interface to these resources is provided by different Grid middlewares.

The remainder of this paper is organized as follows. Related work is reviewed in Section 2. Section 3 presents a general overview of Grid technology and a description of the two Grid middlewares used in this project. Different thoughts on how to solve the problem is presented in Section 4. Section 5 describes different issues related to the benchmarks and the respective Grid middlewares. Benchmarks and results are shown in Section 6. Section 7 discusses the results and future work.

2 Related Work

Different tools have been developed to transparently handle the dynamic nature of Grid systems as well as standards for developing grid integrated applications. Efficient scheduling and execution of PSA's on a Grid is a big challenge for Grid developers. These applications often consist of a large number of jobs where the final result is dependent on all the individual results. These jobs must therefore be scheduled and distributed effectively so not to delay the total execution time. Methods including re-use of common files between executions and adaptive execution to migrate jobs to provide better resources are some suggested solutions [9, 10].

There are not many tools available for benchmarking grid environments. Their heterogeneous nature makes this challenging compared to traditional parallel high performance systems. Liang Peng et al.[14] have done some work on benchmarking the performance between SGE and Globus in terms of CPU utilization and turnaround time. They noticed that the overhead introduced by the grid middleware is negligible for large problem sizes. In their case it actually changed very little even if the problem size grew significantly. For short jobs, however, they found that the overhead can be quite significant, sometimes half the turnaround time for a job. They also found that the Globus middleware generally had more overhead than SGE. The significant overhead for small jobs is what we are trying to minimize in our PSA.

George Tsouloupas and Marios D. Dikaiakos [19] suggest a method for ranking resources in a grid according to a ranking function. They have developed a tool called *SiteRank*, a module built on top of GridBench. With it a user can rank all resources in a grid with respect to a specific application. This tool can be used to better utilize the resources in a grid for any kind of job including the short ones presented in this paper.

3 Grids

Grids are often referred to as High Throughput Computing (HTC) [1]. The Grid is a collection of different, privately owned, computer resources to form a type of heterogeneous, virtual supercomputer [2] for providing computing power for large-scale jobs. Their job is to distribute a high number of jobs efficiently through a common interface and provide long lasting computation time for them.

A simple Grid can be one formed by local workstations, for example inside an office environment. Every day there are hours of idle computer time during for example the lunch hour, staff meetings, after-office hours and at night. During these idle times, the workstations could be used to shorten a Grids job queue. If the local resources are not enough, the number of resources can be dramatically increased by connecting the local grid to remote grids in other locations. The ultimate grid would be the one with access to all the computing resources in the world. However, people are usually very reluctant to let other unknown people use their hardware, at least while they are using it themselves. So what can be done to aid in the process of getting permission to use these resources?

Introducing a computational economy [3] is suggested as a good motivator for people sharing their resources over the Internet making it a computational power grid. This also opens for smaller companies to buy only the resources they need to get their current jobs done and not make huge investments in own infrastructure and computing power.

Other features needed in a Grid to attract users and resource providers are a transparent interface for resource allocation and administration, fault tolerance and different security and authorization tools, so they know their resources will not be exploited [2].

Grids are, however, not to be confused with clusters. A cluster is typically a collection of identical nodes with the same processor and OS, typically containing a static number of nodes, all placed in the same physical location. A Grid is a heterogeneous system [18] with different types of nodes (e.g. computational or storage nodes), processors and OS's. Grids are dynamic in number and resources while clusters are generally more static over time. Another important feature is that Grids can contain different HPC environments such as clusters and supercomputers in addition to other resources [13]. Thus, a Grid can represent a heterogeneous environment with the possibility to utilize the power of supercomputers for less embarrassingly parallel tasks.

Grids can be a cheap alternative to dedicated supercomputers, since a Grid can utilize idle time already available from workstations. These workstations can be very cheap and don't need special rooms or cooling facilities like large supercomputers do (unless a large number of nodes are clustered together). However, as grids are heterogeneous systems, they are best suited

for embarrassingly parallel jobs where the individual jobs are independent of each other during execution. There are several cases where a collection of heterogeneous workstations is not a suitable replacement, e.g. for fine grained parallel tasks with high dependency between processes. Here, each job is best run on its own cluster or supercomputer for optimal performance. Submitting multiple fine grained jobs simultaneously to a grid with access to multiple clusters or supercomputers can give a combination of coarse- and fine grained job execution. For example a large job consisting of multiple independent fine-grained jobs can be automatically distributed by a grid and be run coarse grained on different clusters simultaneously [4].

As an example, at CERN they are currently developing grid tools for their Large Hadron Collider (LHC). They need an incredible amount of storage and computation power and are connecting sites all over the world to their grid to satisfy their need [7]. Without a Grid, it would be impossible to maintain the data and computation throughput necessary for the LHC project.

3.1 Grid features

3.1.1 Condor 6.8.6

Condor is a free grid manager from the Condor team. It was born at the University of Wisconsin in the 1980's as a combination of a doctoral thesis on cooperative processing, the Crystal Multicomputer and Remote Unix. It creates a High-Throughput Computing (HTC) environment by opportunistically utilizing workstations connected through a regular network, remote as well as local. The main features that makes this environment possible is ClassAds, Checkpointing & migrating and Remote System Calls [18, 15, 12].

The ClassAds system is a powerful mechanism for matching jobs to execution nodes. Users advertise their resource needs for a job and Condor matches them with the resource ads for the available workstations. This way the necessary resources are acquired to best match each job.

Checkpointing is a system for transparently moving running jobs from one workstation to another, if necessary. This will for example happen when a user returns from lunch and starts using his or her workstation. Condor will only schedule a job to a node which has been idle for a predefined amount of time, thus not bothering the owner of the respective hardware. Each running job is regularly and transparently checkpointed during execution to make this possible. When a job checkpoints to another node, the new node can resume execution from the last checkpointed state.

Remote System Calls technology, as with checkpointing, requires re-linking of the job executable with specific Condor libraries. The Remote System Calls feature preserves the submitting node's local execution environment, by redirecting a jobs I/O mechanisms back to the submitting node.

Thus, distribution the executables and its input files are not necessary prior to job execution, as this is handled automatically. It also gives a user access to the executing node without having a login account on it. There are however a number of limitations¹ to jobs which are to support checkpointing (including running them on Windows nodes). If for some reason the executable cannot be relinked to run in the standard Condor universe (e.g. no access to source files), the executable can be run unaltered in the “vanilla” universe.

To run jobs with dependencies, Condor includes a feature called DAGMan. This is a Directed Acyclic Graph Manager where rules for job dependencies and pre- and post processing scripts can be set up in a special file. The pre- and post scripts are run locally on the submit host. When this type of job is submitted, the DAGMan takes care of the order of execution according to the rules specified by the user. However, each job defined in a DAG still needs its own regular Condor submit script.

Distributed Resource Management Application API (DRMAA) 1.0 Java and C bindings are also supported. This makes it possible integrate Grid technology into applications, instead of manually submitting jobs through a console.

Many other projects are under development for use with Condor, including a file handling system called Stork, a system monitoring tool called Hawkeye, MW, a master-worker paradigm for Condor and more.

3.1.2 Sun Grid Engine 6 (SGE)

Sun grid engine is a free Grid manager from SUN [5]. It is now an open source project, with support contracts available from Sun. One of the newest features in v6.1 is Resource Quotas, a feature for controlling resources in the grid. Access rules to different parts of a grid can be set up for users, groups, projects, etc. for fine grained control of the available resources. A Condor ClassAds alternative in SGE is Boolean operations. This is a tool for creating rules for specifying resource needs with AND, OR and NOT operations.

Job dependencies can be managed using the Grid Engine Array Job Interdependency (ARI)² feature. This, in combination with prolog- and epilog scripts, gives similar functionality for SGE, as DAGMan does for Condor.

Execution of parallel jobs (MPI or PVM) is supported through a dedicated interface, as with Condor. SGE also supports checkpointing and migration among other tools and functions.

A GUI interface for easy configuration and administration of queues, jobs and nodes is available for Linux. However, all features in this GUI are also available from the command line.

¹http://www.cs.wisc.edu/condor/manual/v6.8.5/1_4Current_Limitations.html

²<http://gridengine.sunsource.net/news/GE61ARIsnapshot-announce.html>

Distributed Resource Management Application API (DRMAA) 1.0 Java and C bindings are supported on the Linux platform, but not on Windows.

4 Implementation Ideas

In this section, different ideas on how the respective Grid middlewares can be used to support our application, is discussed. Different issues concerning the location of input and output files are also considered at the end.

Our PSA may consist of thousands of permutations, where each needs to be computed to find the final result. Luckily, Grid systems support methods for submitting multiple jobs automatically, as described below.

In a Grid environment, all nodes must have access, locally or remotely, to all possible resources needed the tasks they are given. If not, the tasks will obviously return an error. In most cases, these files are located in remote places, e.g. on a NFS server for easy administration, and have to be transferred to the respective nodes for each task. Since the tasks in our application are so short, this extra file transfer overhead is factor to be considered. Since many files are common between the different tasks, methods for collecting multiple tasks in one job is the main focus in this project.

Different methods are considered to minimize overhead:

- Altering the source code to execute multiple calculations from input in the argument list
- Tuning the schedulers for faster submitting of jobs
- Bash script with multiple executions
- DAGMan and Master-Worker features in Condor

4.1 A golden rule

A golden rule is to exploit application domain optimizations before platform domain ones. Altering the source code is not a general solution across applications, but in many cases it should doable to make an application execute multiple tasks by altering the input arguments. By altering the code to include multiple computations, result comparison between tasks can be done internally on each node, and only the best result has to be returned. However, how much computation time is really saved by making the application do multiple executions internally compared to executing multiple single-executions in a bash script? Time will be saved by not having to start and stop the executable for every task. The question is if the time saved for each execution noticeable compared to just specifying multiple runs in a script.

This might come down to which is easier in the long run, depending on the application. Defining multiple runs in a script, each with different input, or alter the source code to make it possible to submit one single executable which does multiple runs internally.

4.2 Scheduler tuning

The default scheduler settings might not be optimal for every compute farm scenario. Different actions can be taken to fine tune the schedulers for optimal performance in specific environments. The SGE scheduler supports different tools³ for debugging and validation of scheduled jobs. These can be turned on or off, depending on the specific needs. When the Grid is in production state, these tools might not be necessary all the time and can be turned off by the administrator.

For example, configuring the SGE scheduler for immediate scheduling, will increase the throughput of the compute farm. The only limitation is the power of the machine hosting the master and scheduler. If this machine is overwhelmed by work, the scheduler can be configured to run jobs only in a fixed schedule interval, which also is the default setting.

In Condor, different parameters can be tuned⁴ in the configuration files for faster scheduling. One of Condor's default behaviors is not to schedule jobs to non-idle nodes. It also preempts and/or suspends jobs, if the current node becomes unavailable due to user interaction. These features can be disabled, if seen fit for the compute farm. Parameters like "JOB_START_DELAY", among others, is a candidate for optimization as well.

4.3 Submitting multiple jobs with one submit file

With condor one can submit multiple jobs in one submit file simply by stating the number of jobs with the "Queue" command, e.g. *Queue 50* for 50 jobs. The job can be identified with the \$(CLUSTER) macro and sub-jobs with the \$(PROCESS) macro. In this case, each sub-job gets a unique identifier from 0-49. Different input parameters can be defined in the submit script by the use of these macros. Using this multi submit method is similar to submitting 50 jobs manually, thus it does not remove any scheduling overhead. It only saves the user time by instantly queuing X number of jobs automatically.

For SGE, the alternative is Array Jobs. There is no need to alter the submit script, as the array job is specified as an argument to the SGE submit-to-queue binary *qsub*. Instead of the \$(PROCESS) macro in Condor, SGE defines a set of environment variables for the array job, to identify the task and the task range. To submit an array job, type the following when submitting a job: **qsub -t 1-10:2 script.submit**. This will queue 5 jobs with step

³<http://docs.sun.com/app/docs/doc/820-0698/enfky?a=view>

⁴http://www.cs.wisc.edu/condor/CondorWeek2007/large_condor_pools.html

2. The tasks will get `SGE_TASK_ID` 1, 3, 5, 7 and 9. `SGE_TASK_ID`, `SGE_TASK_FIRST` and `SGE_TASK_LAST` are environment variables set by SGE for this particular array job. The scheduler process is the same as for Condors Queue X command. Each job is scheduled individually, but time is saved by the automatic queuing of multiple jobs at the same time. The environment variables can be used to automatically select the current input files for the respective tasks.

4.4 Bash script with multiple executions

In the submit script in Condor, one can specify the executable as a bash script or directly as a binary executable. Multiple executions can be specified in a bash script and the binary can be transferred as an input file. This way, the executable is only transferred once and is reused by all the tasks specified in the bash script. Each task's result will be appended to the specified output file if written to stdout. If the application creates any new files, these will also be transferred back to the submitter automatically by Condor. Thus, the execution of multiple tasks does not overwrite any intermediate results.

SGE's submit script is a basically a bash script, where `#$ xxx` defines SGE specific flags or options. Multiple executions can therefore be defined directly with bash arithmetic and submitted as is. SGE does not automatically transfer any input files. Its submit script only invokes the remote host's environment as if it was invoked locally. However, SGE supports prolog/epilog scripts that can perform any necessary processing, including file transfer, before or after job execution. These scripts are run on the execution host and not on the submit host as in Condors DAGMan alternative. For SGE, the executable and input files can be located on NFS for easy administration, or locally in the same path on every node for minimum network traffic.

4.5 Condor specific tools

4.5.1 DAGMan

DAGMan makes it possible and easy to define job dependencies. The jobs in the DAG are regular Condor submit scripts and each job is scheduled individually. Therefore, DAGMan does not help to minimize scheduling of jobs in any way. It can, however, help with post execution result gathering, as described in Section 4.6.2.

4.5.2 MW: Master-Worker

The master-worker paradigm [17, 6] can be very quick for a collection of short jobs. The Condor implementation consists of a set of abstract classes, namely Task, Driver and Worker. The Driver sits below your application

and manages a pool of Workers and set of user defined Tasks. The Workers pick up Tasks, does the user defined work on them and returns result to the Driver. The implementation is specific to each application and will therefore involve altering the existing code, if not implemented during initial development of the respective application.

In the case of the Dagoc software, implementing MW would involve an epic amount of work. Therefore, MW should be a paradigm though of and implemented from the start of development and not after the application is coded in a serially.

MW falls outside the scope of this paper, as we are mainly looking at ways to use Grid tools to optimize execution of our application without massive alteration of the source code.

4.6 Handling the input/output files for Dagoc

The amount of files that are needed for each task as mentioned in 1, might become a factor for the total execution time of our small jobs. Obviously, the job executions would benefit from reusing as much of these files as possible on each node. One solution would be to have all the data and the executables locally on all nodes, but this would be difficult to administer. For easier administration all files could be put on NFS, but this might generate a lot of network traffic, since every calculation has to access it for its input and executable.

By sending a bunch of calculations in each job, we could transfer all common files for each calculation once to the respective execution node and save a lot of traffic. This way, each task in the job would only have to transfer a small amount of files, unique to that job, from the NFS server. All other files would be temporarily available locally on the node for the duration of that particular job.

4.6.1 Condor specific file handling

In Condor, one can transfer input files from the submitting client by defining the input files in the submit file with the option: *transfer_input_files = file1,file2,...*. These files will be copied next to the executable in the execution node's spool directory. However, it does not support transfer of whole directory structures, only lists of specific files. See Appendix ?? for a Condor submit script example with file transfer.

Since the "Vanilla" universe is used (see 4.6.4), the lines *should_transfer_files = YES* and *when_to_transfer_output = ON_EXIT*, is needed to specify that the executable and results are to be transferred between submit host and execute host.

4.6.2 Result gathering with DAGMan

When distributing our application onto a Grid, we need to do post processing after all the calculations are finished to find the best result. This can be set up using job dependencies in DAGMan. Each execution job has its own regular Condor submit file, as shown in B.2.2. A solution was implemented where the last job was a simple bash script submitted with “*universe = local*”. This script collected the results from all the return files and put them all in one single file. The respective input parameters were saved as well. Regular Linux tools such as ‘grep’ was used in the first result collection. Then, a short C++ program was used to extract the best result from the new single result file. The input parameters for the job with the best result was returned. Next, Dagoc was executed locally on the submit node, with the input parameters from the best task, to save all the execution data in the database. All this post processing is done automatically, in sequence, by the last bash script.

4.6.3 Result gathering with SGE

Using the ARI functionality mentioned in Section 3.1.2, one can set up a post job result gathering script similar to the one described above. However, if NFS home directories are not used, each execution host will have to return their results to the submit host for final result extraction. File staging⁵ with epilog scripts can be used for this. File staging has to be enabled by the administrator. The epilog script can then use different variables, set by SGE, to identify which files are to be sent where. An epilog script was set up to transfer the output from SGE’s array tasks back to the submitter (see Appendix C).

4.6.4 Problems with NFS file access

There were some problems encountered while trying to relink the Dagoc executable with the Condor libraries, as described in Appendix B.2.1. Hence, the jobs were submitted in the “Vanilla” universe in the benchmarks. Remote system calls and checkpointing was therefore not available.

Accessing SQLite databases on NFS may induce problems as well. In some perfect, up to date NFS setups it might work, but some setups, including ours, have issues with file locking and databases. As Dagoc uses an SQLite database for data access, this problem was encountered while having this database file remotely. A workaround in Condor is to copy the database file as an input file in the submit script. The other alternative is to have a local copy on each node. This should not be a problem for the calculations in our case, as the intermediate tasks are not writing to the database, only

⁵<http://gridengine.sunsource.net/howto/filestaging/filestaging6.html>

reading. However, administration might become cumbersome and jobs may at times give false results if some nodes are not updated correctly with new database files.

4.7 Other areas of interest

- Grid installation procedures (See Appendix A)
- Job delivery/execution methods (See Appendix B)
- Heterogeneous nodes (Hardware and OS/Linux Distro)

5 Model

In this project, the focus is on small office environments with limited resources. A small 3-node grid is used in the benchmarks. The nodes are basic workstations connected through an Ethernet network. The Linux distribution Fedora 7 is used as operating system on each node. Furthermore, the workstations are quite different, as shown in Table 1.

Methods for exploiting multiple processors on multiple nodes in a serial application are considered. If using multicore nodes in a Grid, one would be able to exploit all the cores available, without altering the source code of the application. Grid schedulers are able to schedule one job per CPU on a node. Thus, the [P4Hyper] machine will be scheduled two jobs simultaneously when used in a Grid, since the HyperThreading technology is interpreted as an extra CPU by the operating system. This might give a slight speedup, although not twofold since the extra CPU detected by the OS is only virtual.

All three nodes in the Grid are assumed to be idle. Condor is configured in testing-mode, to eliminate the waiting time for non idle nodes. This is necessary, because one of the execution nodes is also used as the submit node. The default Condor setup excludes non idle nodes as candidates for jobs. This will make the Condor environment more similar to the SGE's. The nodes are not used for other tasks while benchmarking, for not to bias the results. All tasks used in the benchmarks will be given identical input. The calculated results can then be used to verify that the same calculation is executed in every task.

Job execution is handled differently in the two middlewares, as described in Appendix B.1.1. To support SGE's file locality concept⁶, all the necessary input files, including the binary executable, are assumed available on NFS. Each node can find all specified files using the same paths provided by the SGE's bash script. However, NFS file caching and buffering might bias the result in SGE's favor, as some common files might already be available in the local file buffer on the execution nodes for the following task, reducing file

⁶<http://gridengine.sunsource.net/howto/nfsreduce.html>

transfer. Since we are not using NFS mounted home directories, the results from SGE will be located on the local home directories for the submitting user on each execution node. This may become a significant factor in the benchmark results, since Condor automatically transfers all results back to the submitting node. Hence, two benchmarks will be run for SGE; one with and one without result file transfer back to the submit node. An epilog script is configured, in the SGE queue, to be used by each task. After a job completes, the epilog script is executed and the results are transferred back to the submitter using SCP. Passwordless SSH keys were distributed among the nodes, prior to job execution. Thus, connection authentication is handled automatically, without user interaction.

The SQLite database file needed by the application will be locally available on each node during execution. This is due to the NFS file locking problems described in Section 4.6.4. For the SGE tests, a local copy is situated on each execution node. As for Condor, its regular file handling will be used and the database file is transferred with the job, along with the executable binary or bash script. However, the textual input files for each task will be located in the same NFS location for both SGE and Condor. Thus, the extra file transfers for Condor, compared to SGE, will be the binary executable and the database file from the submitting node to the execution node.

After all results are copied back to the submit host, the final result can be collected using the method described in Section 4.6.2. For Condor, DAGMan can be used. However, DAGMan jobs are not scheduled instantly, but only after about 5 minutes. Thus, the benchmarks for Condor were run by simply submitting the computation job script directly, since this is scheduled instantly. The time taken to find the final result is therefore not included in the numbers for either Condor or SGE. However, the final script is run manually and the extra time used is given.

	[Athlon64]	[P4Hyper]	[Athlon32]
Processor	AMD Athlon64	Intel Pentium 4	AMD Athlon
Extras	64-bit support	HyperThreading	N/A
Speed	3500+	3.0Ghz	2500+
RAM	2GB	1GB	1GB
OS	Fedora 7		
Grid job	Submit/Execute	Master/Submit/Execute	Submit/Execute

Table 1: Workstation specifications

5.1 Test parameter

Walltime is used as the parameter to compare the time taken to execute X very short tasks serially versus distributing them on a grid. We are looking

for speedup in the range of seconds and minutes, not clock cycles, hence the choice of timing parameter granularity. The following tests were run:

- Time used serially on each of the nodes
- Average time for 3 runs of single task jobs on both Condor and SGE
- Average time for 3 runs of multi task jobs on both Condor and SGE
- Average time for 3 runs of multi task jobs on SGE without result file transfer
- Altering the job-task ratio for 1000 jobs

The benchmarks were run 3 times to see if the [P4Hyper] machine would give any significant difference in the total execution time, as well as to verify the results. HyperThreading does not nearly give double the computation power, hence two jobs running simultaneously on the [P4Hyper] machine might use longer time altogether than if executed on two different nodes. A benchmark was run at the end, altering the job-task ratio, to see if there is room for fine-tuning the amount of tasks sent to different nodes.

The actual timings are extracted from log files, capturing the submit time and the end time for the last task in the job. In Condor the user specifies the log file name, in which the submit, start and stop times for each task in a job are recorded. A bash script is used to automatically extract and calculate the time used for the total job (See Appendix C).

SGE has a tool, “*qacct*”, which extracts data about jobs, including wall-time. Data from each job is piped to a file and another bash script is used to extract the timing results for the respective jobs (See Appendix C).

6 Benchmarks and Results

This Section shows the results from computing X small tasks serially compared to distributed on a small three node Grid. The ideas from Section 4.3 and 4.4 were used. A discussion follows of the results from using these methods in both Condor and SGE. The serial results are shown first for comparison. As mentioned in Section 5, finding the final result was to be done manually for practical reasons. The time used to collect, compare and extract the best result from 1000 tasks, was about 3 seconds. This is negligible, when the corresponding calculation time is in the range of over thousand seconds, on SGE and Condor respectively. For 100 tasks, it took less than one second.

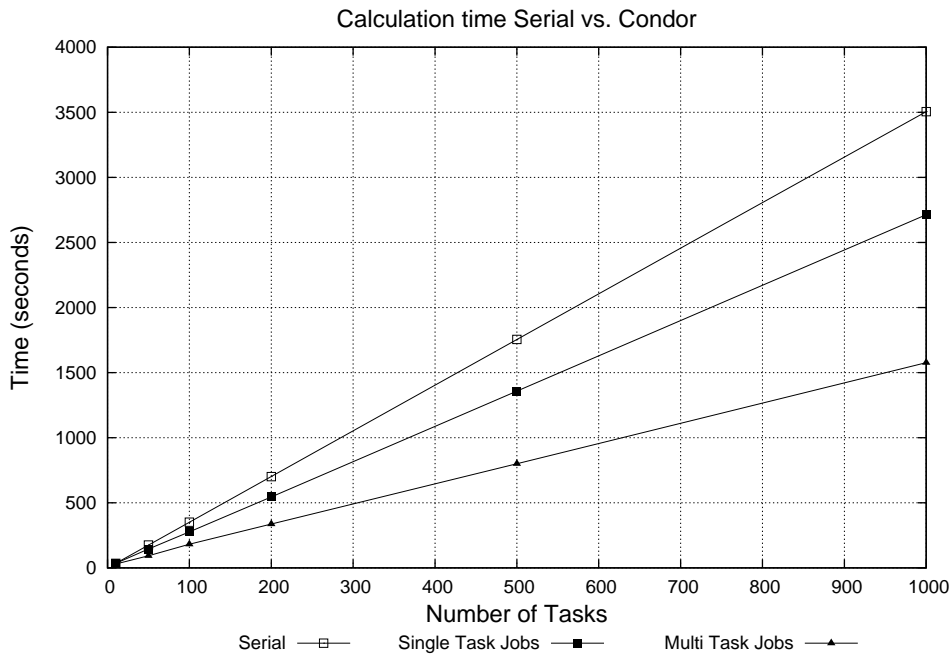


Figure 1: Calculation time serial vs. Condor

6.1 Serial execution

Table 2 shows that the serial execution time varies by about 27% between the fastest and the slowest node. The average time taken between the three was used later in the comparisons. One single task is shown to take between 3 and 4 seconds. The same task was used in all benchmarks and this was verified by checking the result of the tasks.

Machine/nTasks	10	50	100	200	500	1000
Athlon64	30	150	299	599	1498	2988
P4Hyper	37	186	371	741	1853	3702
Athlon32	38	192	382	764	1915	3826
Average(sec.)	35	176	350.67	701.33	1755.33	3505.33

Table 2: Average time for serial execution

6.2 Results for Condor

The first grid benchmark was simply submitting all the tasks in single task jobs, using Condor's *Queue X* command mentioned in 4.3. There is no significant difference in the total time used by the different runs, as shown in Table 3. Thus, the discussion about the HyperThreading capability of one

of the nodes from Section 5.1 may not have any significant implication for single task jobs. However, since all benchmarks are the same, the scheduling will most likely be very similar in each run.

Run/nJobs	10	50	100	200	500	1000
1	34	149	275	543	1353	2713
2	34	142	277	552	1362	2744
3	35	143	278	540	1358	2742
Average(sec)	34.33	144.67	276.67	545	1357.67	2733

Table 3: Condor single task job execution time

Since each task only takes 3-4 seconds, scheduling and file transfer overhead might add a significant delay to the total job execution time. The results in Table 4, compared to 3, show that this is indeed the case. In this test, multiple tasks were sent in fewer jobs, lowering the scheduling and file transfer delays considerably as the task count increased. There is an even smaller difference in time between the different runs in this test than the former. One should think that if the [P4Hyper] machine was given two 100-task jobs while another node is idle would give more difference. Disabling the HyperThreading feature altogether, showed little difference in the timing results for Condor. Furthermore, more benchmarks with different job-task ratios, should be run for a more secure conclusion. This will be considered future work.

Run/Jobs x Tasks	5x2	5x10	10x10	10x20	10x50	10x100
1	30	91	183	338	799	1577
2	30	94	181	336	802	1575
3	30	95	181	335	800	1577
Average(sec)	30	93.33	182.67	336.33	800.33	1576.33

Table 4: Condor multi task job execution time

Table 5 shows the speedup of the two former benchmarks. Sending each task as a single task job peaks at about 1.29 speedup, which is not very good keeping in mind the use of three times the computing power. Multi task jobs, however, show a much higher speedup. Thus, it seems that for such short tasks as in our case (3-4 sec), one can gain a lot from submitting multiple tasks together when using Condor. A speedup of 2.2 is seen for 1000 tasks compared to serial execution, which in turn is a speedup of factor 1.7 compared to single task jobs. Fewer result files are transferred back as well, though their size are larger according to the number of tasks in the respective jobs.

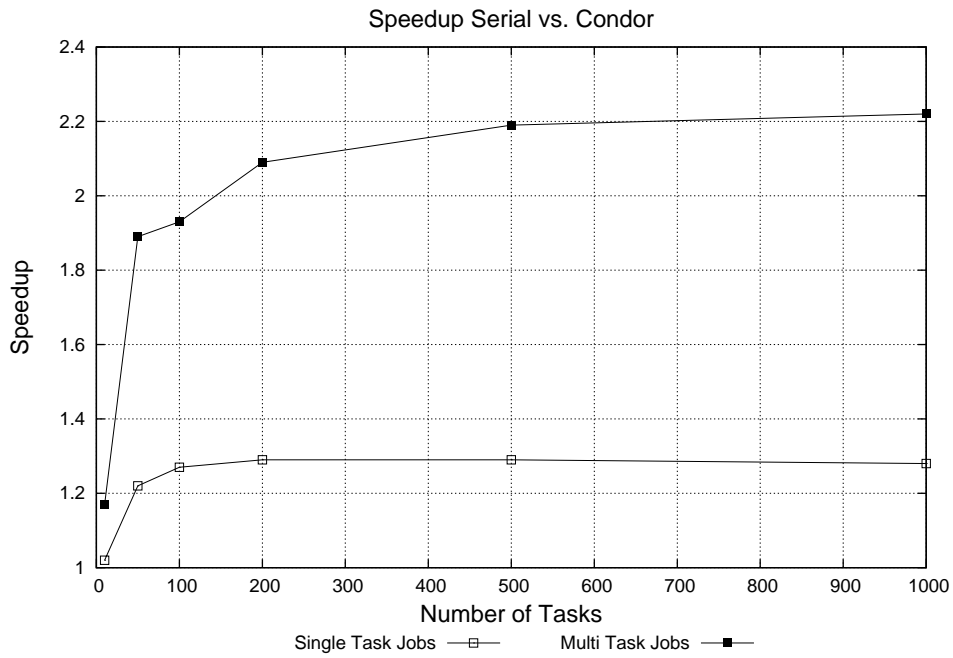


Figure 2: Speedup serial vs. Condor

nTasks	10	50	100	200	500	1000
Single Task Jobs	1.02	1.22	1.27	1.29	1.29	1.28
Multi Task Jobs	1.17	1.89	1.93	2.09	2.19	2.22

Table 5: Speedup using Condor compared to serial execution

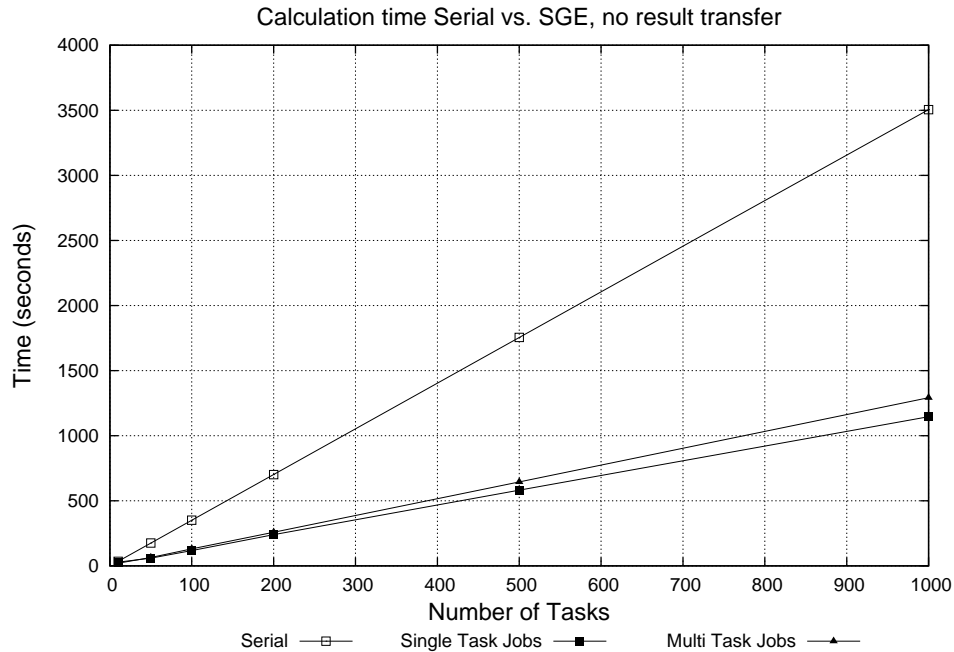


Figure 3: Calculation time serial vs. SGE

Run/nJobs	10	50	100	200	500	1000
1	26	59	117	249	588	1145
2	31	59	116	227	587	1145
3	30	60	118	242	571	1147
Average(sec)	29	59.33	117	239.33	582	1145.67

Table 6: SGE single task job execution time

6.3 Results for Sun Grid Engine (without result file transfer)

Table 6 shows very good results for single task jobs on SGE, even faster than Condor’s multi task job executions. Compared to the average serial calculation time, a speedup factor of the number of nodes used in this small grid is seen. This shows that all the extra file transferring done by Condor creates a significant amount of overhead. Apparently, SGE has nearly nonexistent overhead for these particular tasks when not transferring the results back to the submitter.

Table 7 shows an interesting result. It actually shows slower performance for multi task jobs than for single task jobs. Some tweaking of the job-task ratio was performed and generally showed that the more jobs submitted (with fewer tasks), the closer the timings came to the single task jobs. Thus, it seems that for SGE, single task jobs submitted as array jobs, will perform equal to or better than multi task jobs when not transferring the results back

Run/Jobs x Tasks	5x2	5x10	10x10	10x20	10x50	10x100
1	21	64	130	259	644	1293
2	21	65	130	259	644	1287
3	20	65	131	258	648	1294
Average(sec)	20.67	64.67	130.33	258.67	645.33	1291.33

Table 7: SGE multi task job execution time

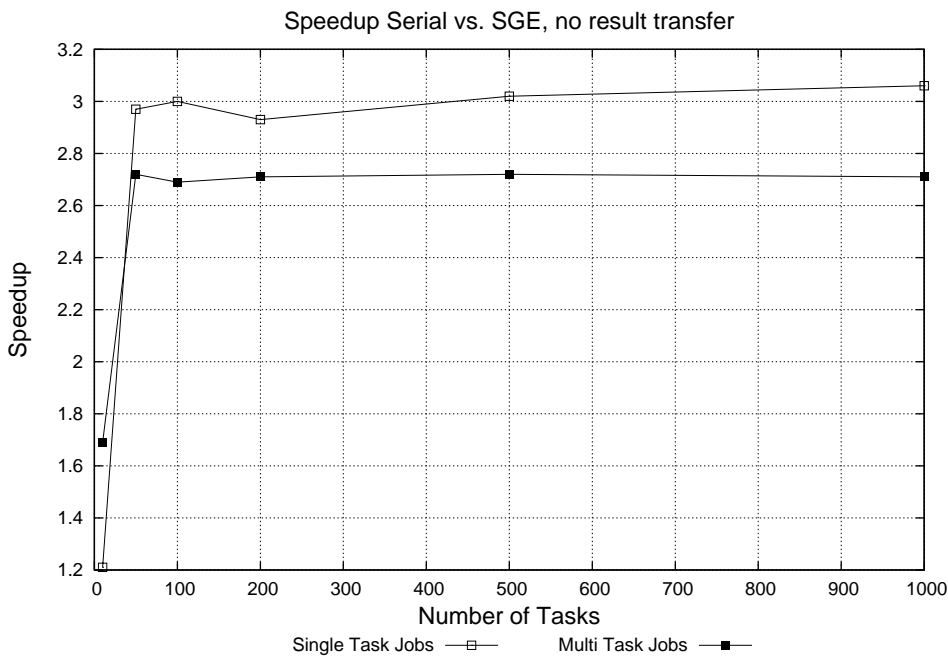


Figure 4: Speedup serial vs. SGE

to the submitter during execution.

6.4 Results Sun Grid Engine (with result file transfer)

In this benchmark, an epilog script is used to transfer the result files for all tasks back to the submit host. This benchmark is run to be fair to Condor, as it transfers all result files back to the submit host automatically. Both stdout and stderr will be transferred for comparison, although one can choose not to transfer the stderr files if they are not needed. Actually, since the epilog transfer script is a regular bash script, one can choose to transfer whatever, in whichever way found suitable.

The results in Table 9, shows the execution time with result transfer back to the submit host. Compared to the single task jobs without result file transfer from Table 6, the numbers are generally higher, especially for small task collections. It is actually slower than serial execution for 10 and

nTasks	10	50	100	200	500	1000
Single Task Jobs	1.21	2.97	3.00	2.93	3.02	3.06
Multi Task Jobs	1.69	2.72	2.69	2.71	2.72	2.71

Table 8: Speedup using SGE compared to serial execution

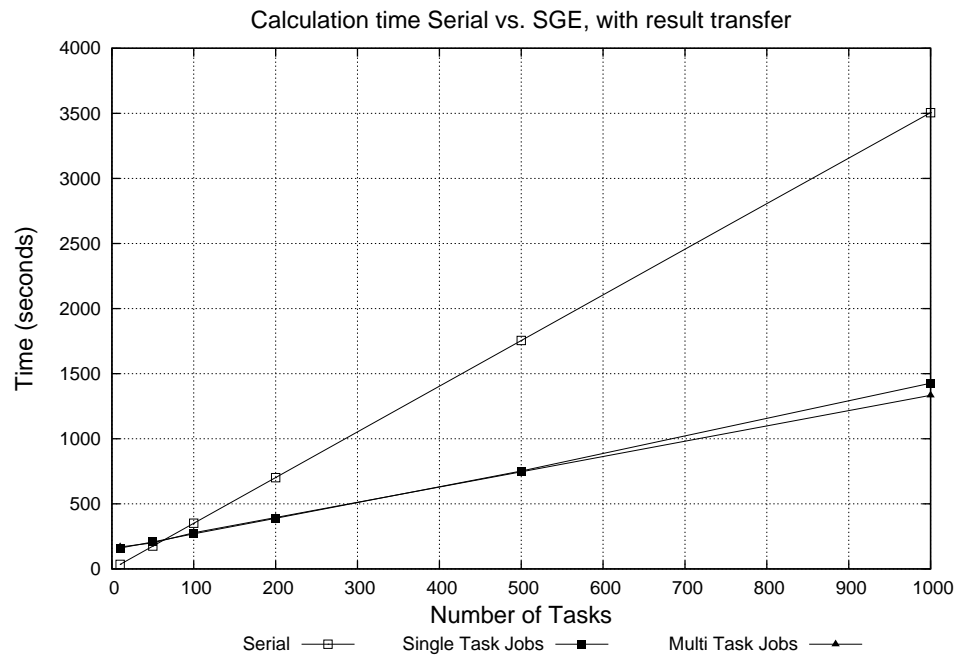


Figure 5: Calculation time serial vs. SGE (with result transfer)

Run/nJobs	10	50	100	200	500	1000
1	160	207	267	388	750	1354
2	161	207	270	389	752	1461
3	161	206	270	390	753	1465
Average(sec)	160.67	206.67	269	389	751.67	1426.67

Table 9: SGE single task job execution time (with result file transfer)

Run/Jobs x Tasks	5x2	5x10	10x10	10x20	10x50	10x100
1	168	200	277	396	746	1334
2	169	200	278	394	746	1335
3	168	200	278	395	745	1332
Average(sec)	168.33	200	277.67	395	745.57	1333.67

Table 10: SGE multi task job execution time (with result transfer)

50 tasks.

In Table 10, it is observed that the multi task jobs with the chosen job-task ratios, perform almost equally well as single task jobs. This result is quite different from the former benchmark, shown in Table 7, where much slower performance was observed for multi task jobs than for single task jobs. Thus, it seems that the extra file transferring levels out the performance between the two.

In Table 11, it is observed that when transferring result files for SGE, the performance is almost identical for both single task and multi task jobs. In Section 6.5, it is observed that this was arbitrary, and that it is possible to tweak the submit scripts in favor of multi task jobs.

6.5 Results when altering the Job-Task ratio

Table 12 shows that there is room for fine tuning the job-task ratio for better performance. It was observed, during the benchmarking, that the [P4Hyper] machine was usually the last machine doing computations, on two jobs simultaneously. Hence, for jobs with large amounts of tasks, the two other machines were idle for a long time. Thus, it seems that Hyper Threading might actually give worse performance altogether, when used on nodes in a Grid. However, this machine could also be looked at as two slower

nTasks	10	50	100	200	500	1000
Single Task Jobs	0.22	0.85	1.30	1.80	2.34	2.46
Multi Task Jobs	0.21	0.88	1.26	1.78	2.35	2.63

Table 11: Speedup using SGE compared to serial execution (with result transfer)

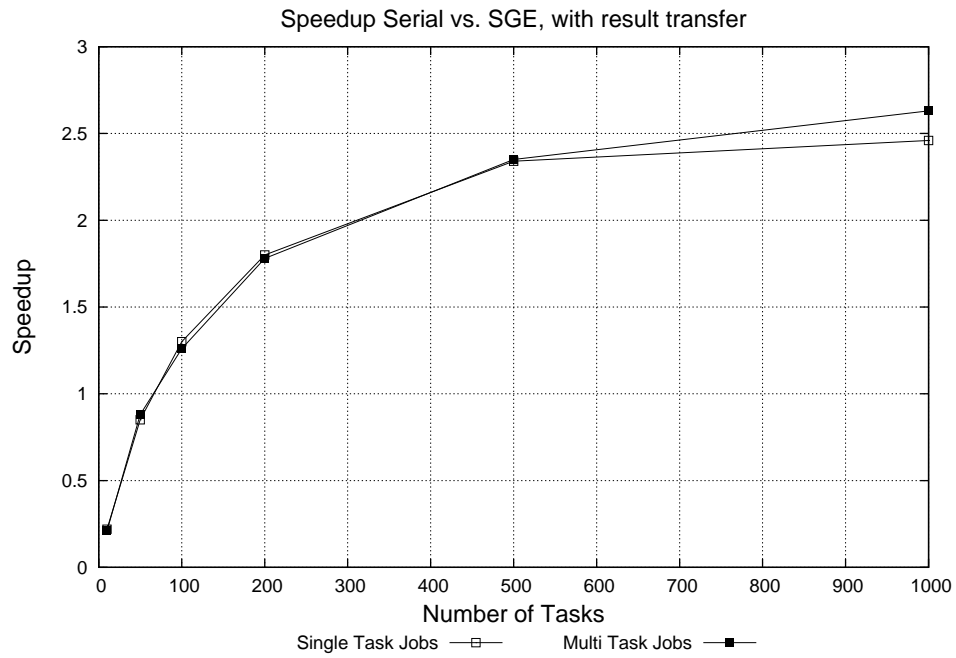


Figure 6: Speedup serial vs. SGE (with result transfer)

Grid/Jobs x Tasks	10x100	20x50	40x25	100x10
Condor	1576	1477	1462	1535
SGE (file transfer)	1333	1290	1246	1263

Table 12: Results for 1000 tasks when altering the job-task ratio

machines, since each of its two jobs take about twice as long to finish as one on the [Athlon64]. Thus, it seems that this environment could benefit from dynamic scheduling, giving larger jobs to more powerful machines.

In any multi user Grid environment, the idle machines would be used to compute tasks from other jobs, by other users. Thus, when speaking of overall throughput in a Grid, this discussion is not equally important. The otherwise idle nodes will be utilized as long as there are other jobs in the queue.

Furthermore, the trend today is multicore processors without Hyper Threading technology, eliminating the discussion altogether. However, the dynamic scheduling idea still stands.

7 Conclusion

In this project, we have seen how a commercial application, developed serially without any initial thought of parallelism or distributed calculation functionality, can benefit from being used in a Grid environment. Two well known Grid middlewares, Condor and Sun Grid Engine, were considered in the process, and relevant functionality was discussed. A discussion of installation procedures and problems can be found in the Appendices as well as practical job submission examples for the respective middlewares.

Due to the short execution time of our tasks, different methods for minimizing scheduler overhead were proposed. In this project, however, only methods using regular Grid submit scripts were analyzed. Thus, all speedup results are gained without altering the source code in any way.

Other methods were proposed, including altering the source code of our application to make it execute multiple tasks internally, tuning the Grid schedulers, and implementing the Master-Worker paradigm.

The first alternative would entail altering the input parameter list and the source code corresponding to the task computations. The alteration would to execute multiple iterations of tasks internally, each with different input data. This is believed not to have considerable speedup compared to our multiple task job proposal. The only time saved is assumed to be the starting and stopping of the executable for each task, and, if internal result comparison is implemented, fewer files to be compared by the last result script. However, this is only an assumption and needs further evaluation before a concrete conclusion can be taken.

Grid schedulers have multiple parameters and features which can be fine-tuned in different ways. By removing unnecessary features and fine tuning different timing constraints, scheduler overhead can be reduced. Removable and tunable features include scheduler monitoring, job validation, load adjustments and different scheduling timings. More information can be found on the web pages for respective Grid systems.

Benchmark/nJobs	10	50	100	200	500	1000
Condor single task	1.02	1.22	1.27	1.29	1.29	1.28
Condor multi task	1.17	1.89	1.93	2.09	2.19	2.22
SGE single transf.	0.22	0.85	1.30	1.80	2.34	2.46
SGE multi transf.	0.21	0.88	1.26	1.78	2.35	2.63
SGE single no-transf.	1.21	2.97	3.00	2.93	3.02	3.06
SGE multi no-transf.	1.69	2.72	2.69	2.71	2.72	2.71

Table 13: Speedup for all benchmarks compared to serial

Implementing the Master-Worker paradigm proposed for Condor, is considered to give easy access to a heterogeneous environment. Implementing it, will entail altering large portions of the source code in our application and is outside the scope of this project. The MW-paradigm describes three classes that have to be implemented, namely Driver, Task and Worker. These classes are used to describe, generate and execute tasks coherently and fast. It is a solution that would have profited from being implemented from the birth of our application, though it is shown to be easily implemented in other serial applications [6], with good results.

In Table 13, the speedup from all benchmarks are listed for convenience. It is observed, that the speedup quickly peaks at around 3, using SGE on three heterogeneous nodes. However, this when looking at the details of this particular benchmark, this speedup does not include the transfer of results from the execution nodes back to the submit node. Taking file transferring into the equation, the speedup, for SGE, is 2.46 for 1000 single task jobs. Another interesting result is actual slowdown when submitting a small number of tasks on SGE. Condor, however, is observed to have much better performance for small number of tasks.

The results show that SGE is about twice as fast as Condor for single task jobs, when transferring the results back to the submitter. Overall, Condor had bad performance when distributing the tasks in our application, using single task jobs. Using three nodes, only a peak speedup of 1.29 was observed. To make it more fair, an extra benchmark was run for Condor at the end for comparison. Using the same bash script as for SGE’s single task jobs, eliminated the transfer of the binary executable and the database file, as opposed to the regular Condor execution. This gave a speedup of 1.43 for 1000 tasks, which is only slightly better compared to the original Condor benchmark. Thus, it seems that transferring the binary from the submit node or fetching it from an NFS server, yield similar results.

In Section 6.5, it is observed that there is room for fine-tuning the job-task ratio for multi task jobs. A performance increase of 7-8% was observed, using 40 jobs with 25 tasks each compared to 10 jobs with 100 tasks each.

These results show that the effort needed for installing a small local Grid

system in an organization, may be well worth it. Automatic distribution of tasks to nodes with idle CPU cycles, gives effective utilization of already available computing resources. For certain applications, no source code needs to be altered to make use of multiple distributed resources.

7.1 Future work

In future work, a dynamic scheduling of the multi task jobs, sized to better match the different nodes in the grid would be of interest. From Table 2 it shows that the Athlon64 machine is about 27% faster than the Athlon32. Should the Athlon64 receive bigger jobs with more computations than the Athlon32 so they finish a job at the same time? Will it level out automatically when enough jobs are in the queue as the Athlon64 most likely will end up receiving more jobs since it finishes them faster?

Methods for dynamic scheduling is found important for heterogeneous environments, like Grids. Different methods are proposed to handle different aspects of heterogeneous environments. Proposed methods include handling dynamic network bandwidth and decreasing makespan of meta-tasks of different size [16, 8, 11]. Combining SiteRank [19] with dynamic size scheduling would be an interesting development.

The benchmarks in this paper did not include the collection of the end results. The user will expect the same end results in the database as for serial execution without having to manually enter it. Without any means of automatically extracting and saving the end result from the distributed calculations, users may become more reluctant towards using Grid technology. The time saved in distributing calculations is lost in collecting and extracting the end result. For Condor, using DAGMan possible solution, where one could add a result gathering job as dependent on the rest of the calculations. For SGE, the newly released ARI functionality could be used to add post jobs dependent on an array of jobs. These methods are only proposed in this project, and not thoroughly tested.

References

- [1] J.H. Abawajy. Job scheduling policy for high throughput computing environments. *Parallel and Distributed Systems, 2002. Proceedings. Ninth International Conference on*, pages 605–610, 17-20 Dec. 2002.
- [2] Mark Baker, Rajkumar Buyya, and Domenico Laforenza. Grids and grid technologies for wide-area distributed computing. *Softw. Pract. Exper.*, 32(15):1437–1466, 2002.
- [3] R. Buyya, D. Abramson, and J. Giddy. Economy driven resource management architecture for computational power grids, 2000.

- [4] EGEE. <http://egee.cesnet.cz/en/info/applications.html>.
- [5] Sun Grid Engine. <http://www.sun.com/software/gridware/>.
- [6] Goux, J.-P., Kulkarni, S., Linderoth, J., and Yoder, M. An enabling framework for master-worker applications on the computational grid. In *The Ninth International Symposium on High-Performance Distributed Computing*, pages 43–50, May 2000.
- [7] GridCafé. <http://gridcafe.web.cern.ch/gridcafe/gridatcern/lcg.html>.
- [8] B. Hamidzadeh, D. J. Lilja, and Y. Atif. Dynamic scheduling techniques for heterogeneous computing systems. *Concurrency: Practice and Experience*, 7(7):633–652, 1995.
- [9] J. Herrera, E. Huedo, R. S. Montero, and I. M. Llorente. Execution of typical scientific applications on globus-based grids. In *ISPDC '04: Proceedings of the Third International Symposium on Parallel and Distributed Computing/Third International Workshop on Algorithms, Models and Tools for Parallel Computing on Heterogeneous Networks (ISPDC/HeteroPar'04)*, pages 177–183, Washington, DC, USA, 2004. IEEE Computer Society.
- [10] Eduardo Huedo, Ruben S. Montero, and Ignacio M. Llorente. Experiences on adaptive grid scheduling of parameter sweep applications. *pdp*, 00:28, 2004.
- [11] Hung-Yuan; Liu Kang-Yuan; Chang Gei-Ming; Lien Chin-Chih Lee, Liang-Teh; Chang. A dynamic scheduling algorithm in heterogeneous computing environments. *Communications and Information Technologies, 2006. ISCIT '06. International Symposium on*, pages 313–318, Oct. 18 2006-Sept. 20 2006.
- [12] Miron Livny, Jim Basney, Rajesh Raman, and Todd Tannenbaum. Mechanisms for high throughput computing. *SPEEDUP Journal*, 11(1), June 1997.
- [13] Liang Peng, Lip Kian Ng, and Simon See. Yellowriver: A flexible high performance cluster computing service for grid. In *HPCASIA '05: Proceedings of the Eighth International Conference on High-Performance Computing in Asia-Pacific Region*, page 553, Washington, DC, USA, 2005. IEEE Computer Society.
- [14] Liang Peng, Simon See, Jie Song, Appie Stoelwinder, and Hoon Kang Neo. Benchmark performance on cluster grid with ngb. *ipdps*, 18, 2004.
- [15] Condor Project. <http://www.cs.wisc.edu/condor/>.

- [16] Prashanth C SaiRanga and Sanjeev Baskiyar. A low complexity algorithm for dynamic scheduling of independent tasks onto heterogeneous computing systems. In *ACM-SE 43: Proceedings of the 43rd annual Southeast regional conference*, pages 63–68, New York, NY, USA, 2005. ACM.
- [17] MW Team. <http://www.cs.wisc.edu/condor/mw/>.
- [18] Douglas Thain, Todd Tannenbaum, and Miron Livny. Distributed computing in practice: the condor experience. *Concurrency - Practice and Experience*, 17(2-4):323–356, 2005.
- [19] George Tsouloupas and Marios D. Dikaiakos. Grid resource ranking using low-level performance measurements. In Anne-Marie Kermarrec, Luc Bougé, and Thierry Priol, editors, *Euro-Par*, volume 4641 of *Lecture Notes in Computer Science*, pages 467–476. Springer, 2007.

A Setting up the grid

A.1 Sun Grid Engine 6.1

SGE master host and execution host was successfully installed on Fedora 7 by following the installation manual found on the Sun web page [5]. After working out the problems mentioned in A.1.3, the installation was pretty straight forward. There were a lot of steps to set up different things and it is a good idea to have a plan or basic idea of the grid before installing.

To install an execution host, it is necessary to copy all the installation files from the master host to the execution host after the master host installation. This way the execution host will get the correct settings set up for the master host.

A.1.1 Possible node configurations

There are several different node functions available for SGE:

- Master host: this is where the grid engine runs.
- Shadow master host: this is a backup host if the master host fails. There can be several shadow hosts in a grid.
- Administration host: nodes that can do administrative tasks on the grid.
- Submit host: a node which can submit and control jobs.
- Execute host: a node where jobs are executed.

A.1.2 Access to files

Each node in a SGE grid either needs all executables and input files locally or as a NFS/AFS mount. There is no automatic way for SGE to copy the executables and input files to the executing nodes. Pre and post scripts can be defined for each job which can transfer files or do other operations needed by the job before and after the execution.

A.1.3 Problems

Some problems were encountered while installing the master host:

- SGE needs the libXm.so.3 library, which can be found in the OpenMotif package, for its GUI application. OpenMotif-2.3.0.0.fc7.ccrma.i386.rpm for fedora core 7 was installed which had the newer version, libXm.so.4, of the library. I had to make a symlink to this file for the SGE GUI to work (since it is a newer version the linker is happy):
ln -s /opt/openmotif/usr/lib/libXm.so.4.0.0 /usr/lib/libXm.so.3
- After unpacking the files, the command **setfileperm.sh \$SGE_ROOT** is to be run to set the right permissions. This failed because of a wrong GLIBC version in Fedora 7. To fix this, open the file "*\$SGE_ROOT/util/arch*" and edit line 248 from 3|4|5) to 3|4|5|*) and run the script again. (NB! This problem did not appear on a machine running Kubuntu)

A.1.4 Windows restrictions

Windows machines cannot run as master hosts, shadow master hosts or scheduler. Windows is therefore limited to execution and submit hosts. Certificates (Certificate Security Protocol (CSP)) are also necessary for communication between master host and windows execution host. The GUI tool *qmon* and DRMAA is not supported either.

A.2 Condor-6.8.6

I installed Condor with a rpm package on Fedora 7, with a tar.gz package on Kubuntu and with a MSI package on WindowsXP. The only prerequisite on Linux was an older version of the libstdc++ library. This can be installed with **yum install compat-libstdc++-33** for Fedora 7 or **apt-get install libstdc++5** on Kubuntu. The installation was easily completed by following the installation manual. The installation of a 3 node Condor pool (with a WindowsXP node) ready to submit jobs was an hours process, where on SGE it was more a days process for a Grid beginner.

Registering an execute and submit node on the master host in Linux is easy, just enter the following command on the host to add:

```
./condor_configure --central-manager=host@domain.com --type=execute,submit
```

Make sure all the nodes networking is set up correctly before running this command. If not, the request might not be detected by the master host.

A.2.1 Possible node configurations

- Master host (even Windows nodes, unlike SGE)
- Execute host
- Submit host
- Or any combination of these

A.2.2 Access to files

In Condors standard universe, access to files, input and output, is handled automatically through remote system calls. In other universes, vanilla, Java and MPI, access is presumed, as default, to be through a shared file system on UNIX machines. If no shared file system is available, file transfer has to be specified in the submit files by the user. Add the following lines to enable file transfer (other options are available):

```
should_transfer_files = YES
when_to_transfer_output = ON_EXIT
transfer_input_files = file1 , file2 , ...
```

A.2.3 Heterogeneous nodes

In a heterogeneous Grid it is beneficial to have as much information about each node as possible. If a job has specific needs it special care should be take to which node the jobs is sent to. Condor solves this with the ClassAds system. Each node has a set of parameters (e.g. Total Memory, OpSys, Architecture, Disk Space) which can be used as requirements for jobs. Condor administrators can specify their own ClassAds for each node as well. Inserting the following in a submit file states that the job needs the Linux opsys and flavor RedHat9 (*OPSYS_FLAVOR="RedHat9" has to be defined in the execution nodes configuration file for it to be eligible*):

```
REQUIREMENTS = (OpSys == "LINUX") && (OPSYS_FLAVOR =?= "RedHat9")
```

This is useful if the executable is compiled only for a specific flavor of Linux and could possibly fail if executed, for example, on Debian. A neat trick to be able to use more nodes would be to have different executables for

different flavors and specify which executable should be transferred to the different nodes. The following line will choose the correct executable for the specific requirement.

```
Executable = exec.$$ (OpSys).$$ (OPSYS_FLAVOR)
```

The executables must have names *exec.LINUX.RedHat9* or *exec.LINUX.Debian* or whatever other flavor registered in the nodes.

A.2.4 Problems

I installed a net install of Debian on a 3rd node and ran into a problem with ip-addresses and hostnames. The installation set the IP address for the nodes hostname in the */etc/hosts* file to a local one (127.0.1.1). This caused the master to block access for the node, because it had an unknown domain address. Commenting out this line and letting DHCP take care of hostname and corresponding ip-addresses fixed the problem.

A.2.5 Windows restrictions

It is not possible to `condor_compile` Windows applications. As a result, remote system calls and check pointing is not available on this platform. Therefore Windows jobs has to run in the "vanilla" universe. The following has to be added to the submit script to run on windows machines:

```
universe = vanilla
requirements = (OpSys == WINNT50)
```

B Using the grid

B.1 Sun Grid Engine 6.1

To start the SGE daemons, enter the following commands as root:

```
$$SGE_ROOT/name-of-cell/common/sgemaster start
$$SGE_ROOT/name-of-cell/common/sgeexecd start
```

B.1.1 Submitting Jobs

Submitting jobs to SGE is done by sending batch scripts to the master server with the command "**qsub /path/to/script.sh**" or through the GUI interface *qmon*. A computer is not allowed to submit jobs unless it is registered as a "Submit Host". Once this is done, jobs sent from this node will be accepted and put in a job queue.

The results are handled different than in Condor. Condor copies all results back to the same folder on the node it was submitted from. SGE

copies the results to the owners home directory. If the home directories are not NFS mounted the results are put locally on the execution node. This can be tricky, as you don't know which node your job is executed on. Thus, SGE assumes the users home directories are NFS mounted.

SGE does not have automatic job executable transfer as Condor does. Job executable must therefore be available on every node in the path given in the bash script. The easiest way to make sure this is available is to have all the files available on AFS or NFS. However, it is possible to manually define transfer of executables in job pre- and post scripts as in Condor's DAGMan if necessary.

Example of a simple SGE submit script:

```
#!/bin/sh
#
# This is a simple example of a SGE batch script
# Request Bourne shell as shell for job
#$ -S /bin/sh
#
/path/to/executable arg1 arg2
```

B.2 Condor-6.8.6

To start Condor, enter the following commands as root:

```
$CONDOR_ROOT/sbin/condor_master
```

B.2.1 Preparing jobs for execution

To use remote system calls and checkpointing/migration in Condor, the executable needs to be relinked with the Condor libraries. Here a problem was encountered on both Fedora 7 and Kubuntu: "ERROR: Internal ld was not invoked! Executable may not be linked properly for Condor!". A solution was not found and the jobs were run in the "vanilla" universe instead.

B.2.2 Submitting Jobs

Submitting jobs is done with the command "**condor_submit job.cmd**". The .cmd file contains different job settings as input/output file locations, files that have to be transferred to the execution node, the required architecture for running a binary etc. This way a user knows that the binary and input files are provided with the job. Copying directories of input files is not supported in the current version of Condor. For jobs with directories as input data a shared file system, AFS or NFS, can be used for input files instead. ClassAds can be used to find a node which has access to the specific remote location.

Example of a simple Condor submit script that copies the executable to the execute node and returns the results to the submit node:

```
#  
# Test Condor command file  
#  
universe = standard  
executable = name_of_executable  
output = executable.out  
error = executable.err  
log = executable.log  
arguments = arg1 arg3  
queue
```

See manual for more options:

http://www.cs.wisc.edu/condor/manual/v6.8/condor_submit.html

B.2.3 Different settings

Condor can either be set up to run jobs on any node idle or not, or to nodes which have been idle for more than 15 minutes (either no keyboard or mouse movement or CPU idle). When no longer idle, the job can be checkpointed and kept on the same node until idle again, or the job can be sent to another idle node. During testing it is recommended to disable the idle settings as was done while benchmarking in this paper.

C Scripts

Here follows a listing of all scripts used in benchmarking Condor and SGE. The last listing is the C++ program used to extract the best result from the pre formatted result file, described in Section 4.6.2.

Algorithm 1 Single Task Job script for SGE

```
#####
# Single Task Job Script for Sun Grid Engine
# Single job submission:
# qsub this_script.sh
# Array job submission:
# qsub -t [t_first]-[t_last]:[t_stepsize] this_script.sh
#####
#!/bin/sh
# Request Bourne shell as shell for job
#$ -S /bin/sh
v1=$(date +%s)
# Run dagoc
for ((i=0; i<[set_num_tasks_here]; i+=1)) ; do
/mnt/project/dagoc -c -start=1 -stop=10 \
"/mnt/project/setups/TEST_sge.sup"
done
# Print Job Data
echo start=$SGE_TASK_FIRST stop=$SGE_TASK_LAST \
step=$SGE_TASK_STEPSIZE id=$SGE_TASK_ID
# Print time taken
v2=$(date +%s)
let v3=$v2-$v1
echo "Seconds used for this task: " $v3
```

Algorithm 2 Restult transfer epilog script for SGE

```
#####
# result_transfer_epilog.sh
# Transfers SGE array results to host
#####
#!/bin/bash
echo SGE_HOST: $SGE_O_HOST
echo HOSTNAME: $(hostname -s)
if [ "$(hostname -s)" != "$SGE_O_HOST" ];
then
    echo "Transferring result to host..."
    scp $SGE_O_WORKDIR/dagoc.sh.o$JOB_ID.$SGE_TASK_ID \
        $SGE_O_WORKDIR/dagoc.sh.e$JOB_ID.$SGE_TASK_ID \
        $SGE_O_HOST:$SGE_O_WORKDIR
    rm $SGE_O_WORKDIR/dagoc.sh.o$JOB_ID.$SGE_TASK_ID
    rm $SGE_O_WORKDIR/dagoc.sh.e$JOB_ID.$SGE_TASK_ID
fi
```

Algorithm 3 Script for extracting runtime for SGE

```
#####
# Get time used from Sun Grid Engine
# Usage:
# qacct -j [job_id] > sgejobsummary.txt
# ./this_script.sh sgejobsummary.txt
#####
#!/bin/bash
#Get job submit time
t0=$(grep "qsub_time" $1 | head -n 1 | \
grep -o [0-9][0-9]:[0-9][0-9]:[0-9][0-9])
s0=$(date --date=$t0 +%s)
#Get last job end time
t1=$(grep "end_time" $1 | \
grep -o [0-9][0-9]:[0-9][0-9]:[0-9][0-9] | \
sort -r | head -n 1)
s1=$(date --date=$t1 +%s)
#Get total time used for array job
let s=$s1-$s0
echo "File:" $1
echo "Time:" $s "seconds"
```

Algorithm 4 Condor single task jobs submit script

```
#####
# dagoc.condor
# 100 Single Task Jobs Condor command file
#####
universe      = vanilla
executable    = dagoc
output        = dagoc.out.$(CLUSTER).$(PROCESS)
error         = dagoc.err.$(CLUSTER).$(PROCESS)
log           = dagoc.log.$(CLUSTER)
REQUIREMENTS = (OpSys == "LINUX") && (OPSYS_FLAVOUR != "FC7")
should_transfer_files = YES
when_to_transfer_output = ON_EXIT
transfer_input_files = /mnt/project/dbases/TEST.db
arguments = -c -start=1 -stop=10 /mnt/project/setups/TEST.sup
queue 100
```

Algorithm 5 Condor multi task job submit script

```
#####
# 10 Multi(10x)Task Job Condor command file
#####
universe      = vanilla
executable    = multiTaskJob.sh
output        = dagoc.out.$(CLUSTER).$(PROCESS)
error         = dagoc.err.$(CLUSTER).$(PROCESS)
log           = dagoc.log.$(CLUSTER)
REQUIREMENTS = (OpSys == "LINUX") && (OPSYS_FLAVOUR != "FC7")
should_transfer_files = YES
when_to_transfer_output = ON_EXIT
transfer_input_files = dagoc, /mnt/project/dbases/TEST.db
arguments = 10
queue 10
```

Algorithm 6 Condor multi task job bash script

```
#####
# MultiTaskJob .sh
#####
#!/bin/bash
v1=$(date +%s)
echo "Start: " `date`
for ((i=0; i<$1; i+=1)) ; do
    ./dagoc -c -start=1 -stop=10 /mnt/project/setups/TEST.sup
done
v2=$(date +%s)
let v3=$v2-$v1
echo "Finished: " `date`
echo "Seconds used: " $v3
```

Algorithm 7 Script for extracting runtime for Condor

```
#####
# Get time used from Condor log file
# Usage:
# ./this_script.sh nameOfLogFile.log
#####
#!/bin/bash
# Get submit time and last job end time
t0=$(grep "Job submitted" $1 | head -n 1 | \
grep -o [0-9][0-9]:[0-9][0-9]:[0-9][0-9])
t1=$(grep "Job terminated" $1 | tail -n 1 | \
grep -o [0-9][0-9]:[0-9][0-9]:[0-9][0-9])
# Convert to seconds and find total time
s0=$(date --date=$t0 +%s)
s1=$(date --date=$t1 +%s)
let s=$s1-$s0
# Print results
echo "File:" $1
echo "Time:" $s "seconds"
```

Algorithm 8 Condor DAGMan script

```
#####
# DAGMan script for
# post processing
#####
Job A dagoc.condor
Job B post.condor
Parent A CHILD B
```

Algorithm 9 post.condor script

```
#####
# post.condor
# Post processing job for condor
#####
universe    = local
executable  = post.sh
output      = post.out.$(CLUSTER)
error       = post.err.$(CLUSTER)
log         = dagoc_dag.log
arguments   = 1 10 GROUP.FIELD.FGasa
queue
```

Algorithm 10 The result extraction bash script for single task jobs

```
#####  
# post.sh  
# Post processing script  
# (single task jobs)  
# Arguments:  
# 1: start 2: stop  
# 3: result_tag 4: res_file_name  
#####  
#!/bin/bash  
echo Collecting results from output files  
s0=$(date +%s)  
rm result.txt  
for i in $( ls $4*); do  
    echo JobID: $i >> result.txt  
    grep Setup $i >> result.txt  
    grep $3 $i >> result.txt  
done  
echo -----  
echo Getting best result from collection  
echo -----  
./findBestResult $1 $2 $3  
echo -----  
echo Running dagoc with the best setup  
chmod 775 postDagmanScript.sh  
./postDagmanScript.sh  
rm postDagmanScript.sh  
echo -----  
echo Done!  
s1=$(date +%s)  
let s=$s1-$s0  
echo "Time used for post processing: " $s "seconds"
```

```

/*#####
* findBestResult.cpp
*#####
#include <sstream>
#include <iostream>
#include <fstream>
#include <string>
using namespace std;
//Input args: start stop resultTag
int main (int argc, char **argv) {
    if(argc != 4)
    {
        fprintf(stderr, "Wrong_#_of_arguments!\n");
        return 1;
    }
    //Get input parameters
    int start = atoi(argv[1]);
    int stop = atoi(argv[2]);
    string tagResult = argv[3];

    string tagSetupFile("Setup_file:");
    string line;
    double bestResult = -1.0;
    string setupFilePath;
    ifstream myfile ("result.txt");
    if (myfile.is_open())
    {
        while (! myfile.eof() )
        {
            getline (myfile, line);
            if(line.find(tagSetupFile,0) != string::npos)
            {
                //Save temporary best setup file path
                string tmpSetupFilePath = line.substr(tagSetupFile.
                    length());

                //Read next line if "Setup Path" tag found
                getline(myfile, line);
                if(line.find(tagResult,0) == string::npos)
                {
                    cout << "Did_not_find_result_tag_for" <<
                        tmpSetupFilePath << endl;
                    continue;
                }
            }
        }
    }
}

```



```

//If "Result Tag" found parse the result
    double tmpRes = -1.0;
    string res = line.substr(line.find(":",0)+1);
    istream i(res);
    if (!(i >> tmpRes))
    {
        cout << "Res:_ " << res << endl;
        cout << "Could_not_parse_result " << endl;
        continue;
    }

    //Save new best result and setup file path
    if( tmpRes > bestResult)
    {
        bestResult = tmpRes;
        //Remove last point
        setupFilePath = tmpSetupFilePath.substr(0,
            tmpSetupFilePath.length()-1);
    }
}
}
myfile.close();

//Create script with final dagoc execution
ofstream outputFile;
outputFile.open("postDagmanScript.sh");
if(outputFile.is_open())
{
    outputFile << "#!/bin/bash\n";
    outputFile << "./dagoc_c_start=" << start << "_-
        stop=" << stop << "_ " << setupFilePath << "_>_
        DAGOCRESULT.RES\n";
    outputFile.close();
}
else cout << "Unable_to_open_output_file " << endl;

//Print results
cout << "Best_result:_ " << bestResult << endl;
cout << "Setup_Path:_ " << setupFilePath << endl;
}
else cout << "Unable_to_open_input_file " << endl;

return 0;
}

```