

Synthetic Grid Workloads with Ibis, KOALA, and GrenchMark

Alexandru Iosup¹, Jason Maassen², Rob van Nieuwpoort², and
Dick H.J. Epema¹

¹ Faculty of Electrical Engineering, Mathematics, and Computer Science
Delft University of Technology, The Netherlands,
{A.Iosup,D.H.J.Epema}@ewi.tudelft.nl

² Department of Computer Science,
Vrije Universiteit, Amsterdam, The Netherlands
{Jason,Rob}@cs.vu.nl

Abstract. Grid computing is becoming the natural way to aggregate and share large sets of heterogeneous resources. However, grid development and acceptance hinge on proving that grids reliably support real applications. A step in this direction is to combine several grid components into a demonstration and testing framework. This paper presents such an integration effort, in which three research prototypes, namely a grid application development toolkit (Ibis), a grid scheduler capable of co-allocating resources (KOALA), and a synthetic grid workload generator (GRENCHMARK), are used to generate and run workloads comprising well-established and new grid applications on our DAS multi-cluster testbed.

Keywords: Grid, performance evaluation, synthetic workloads.

1 Introduction

Grid computing's long term goal is to become the natural way to share heterogeneous resources, and to aggregate them into virtual platforms, used by multiple organizations and independent users. With the grid infrastructure starting to meet the requirements of such an ambitious goal [2], the current evolution of grids hinges on proving that it can run real applications, from traditional sequential and parallel applications to new, grid-only, applications. As a consequence, there is a clear need for generating and running workloads comprising grid applications for demonstration and testing purposes.

A significant number of projects have tried to tackle this problem from different angles: attempting to produce a representative set of grid applications like the NAS Grid Benchmarks [8], creating synthetic applications that can assess the status of grid services like the GRASP project [4], and creating tools for launching benchmarks and reporting results like the GridBench project [15].

This work addresses the problem of generating and running synthetic grid workloads, by integrating the results of three research projects coming from

CoreGRID partners, namely the grid application development toolkit Ibis [16], the grid scheduler KOALA [12], and the synthetic grid workload generator and submitter GRECHMARK. Ibis is being developed at VU Amsterdam³ and provides a set of generic Java-based grid applications. KOALA is being developed at TU Delft⁴ and allows running generic grid applications. Finally, GRECHMARK is being developed at TU Delft⁵ and is able to generate workloads comprising typical grid applications, and to submit them to arbitrary grid environments.

2 A Case for Synthetic Grid Workloads

There are three ways of evaluating the performance of a grid system: analytical modeling, simulation, and experimental testing. This section presents the benefits and drawbacks of each of the three, and argues for evaluating the performance of grid systems using synthetic workloads, one of the two possible approaches for experimental testing.

2.1 Analytical Modeling and Simulations

Analytical modeling is a traditional method for gaining insights into the performance of computing systems. Analytical modeling may simplify *what-if* analysis, for changes in the system, in the middleware, or in the applications. However, the sheer size of grids and their heterogeneity make realistic analytical modeling hardly tractable.

Simulations may handle complex situations, sometimes very close to the real system. Furthermore, simulations allow the *replay* of real situations, greatly facilitating the discovery of appropriate solutions. However, simulated system size and diversity raises questions on the representativeness of simulating grids. Moreover, nondeterminism and other forms of hidden dynamic behavior of grids make the simulation approach even less suitable.

2.2 Experimental Testing

There are two ways to experimentally assess the performance of grid systems: *benchmarking* and *using synthetic grid workloads*. Note that currently existing grids prevent the use of traces of real grid workloads: the infrastructure changes too fast, leading to incompatible resource requests when re-running old traces.

Benchmarking is typically used to understand the quantitative aspects of running grid applications and to make results readily available for comparison. Benchmarks comprise a set applications representative for a class of systems, and a set of rules for running the applications as a synthetic system workload. Therefore, a benchmark is a single instance of a synthetic workload.

³ Ibis is available from <http://www.cs.vu.nl/ibis/>.

⁴ KOALA is available from <http://www.st.ewi.tudelft.nl/koala/>.

⁵ GRECHMARK is available from <http://grenchmark.st.ewi.tudelft.nl/>.

Benchmarks present severe limitations, when compared to synthetic grid workloads generation. They have to be developed under the auspices of an important number of (typically competing) entities, and can only include well-studied applications. Putting aside the considerable amounts of time and resources needed for these tasks, the main problem is that grid applications are starting to develop just now, typically at the same time with the infrastructure [13], thus limiting the availability of truly representative applications for inclusion in standard benchmarks. Other limitations in using benchmarks for more than raw performance evaluation are:

- Benchmarking results are valid only for workloads truly represented by the benchmark’s set of applications; moreover, the number of applications typically included in benchmarks [8, 15] is typically small, limiting even more the scope of benchmarks;
- Benchmarks include mixes of applications representative at a certain moment of time, and are notoriously resistant to include new applications; thus, benchmarks cannot respond to the changing requirements of developing infrastructures, such as grids;
- Benchmarks make difficult either the evaluation of one particular system characteristic (high-level benchmarks), or the evaluation of a mix of characteristics (low-level benchmarks);

An extensible framework for *generating and submitting synthetic grid workloads* uses applications representative for today’s grids, and fosters the addition of future grid applications. This approach can help overcome the aforementioned limitations of benchmarks. First, it offers better flexibility in choosing the starting applications set, when compared to benchmarks. Second, applications can be included in generated workloads, even when they are in a debug or test phase. Third, the workload generation can be easily parameterized, to allow for the evaluation of one or a mix of system characteristics.

2.3 Grid Applications Types

From the point of view of a grid scheduler, we identify two types of applications that can run in grids, and may be therefore included in synthetic grid workloads.

Unitary applications This category includes single, unitary, applications. At most the job programming model must be taken into account when running in grids (e.g., launching a name server before launching an Ibis job). Typical examples include sequential and parallel (e.g., MPI, Java RMI, Ibis) applications.

Composite applications This category includes applications composed of several unitary or composite applications. The grid scheduler needs to take into account issues like task inter-dependencies, advanced reservation and extended fault-tolerance, besides the components’ job programming model. Typical examples include parameter sweeps, chains of tasks, DAG-based applications, and even generic graphs.

2.4 Purposes of Synthetic Grid Workloads

We distinguish three reasons for using synthetic grid workloads.

System design and procurement Grid architectures offer many alternatives to their designers, in the form of hardware, of operating software, of middleware (e.g., a large variety of schedulers), and of software libraries. When a new system is replacing an old one, running a synthetic workload can show whether the new configuration performs according to the expectations, before the system becomes available to users. The same procedure may be used for assessing the performance of various systems, in the selection phase of the procurement process.

Functionality testing and system tuning Due to the inherent heterogeneity of the grids, complicated tasks may fail in various ways, for example due to misconfiguration or unavailability of required grid middleware. Running synthetic workloads, which use the middleware in ways similar to the real application, helps testing the functionality of the grids and detecting many of the existing problems.

Performance testing of grid applications With grid applications being more and more oriented towards services [10] or components [9], early performance testing is not only possible, but also required. The production cycle of traditional parallel and distributed applications must include early testing and profiling. These requirements can be satisfied with a synthetic workload generator and submitter.

3 An Extensible Framework for Grid Synthetic Workloads

This section presents an extensible framework for generating and submitting synthetic grid workloads. The first implementation of the framework integrates two research prototypes, namely a grid application development toolkit (Ibis), and a synthetic grid workload generator (GRENCHMARK).

3.1 Ibis: Grid Applications

Ibis is a grid programming environment offering the user efficient execution and communication [5], and the flexibility to run on dynamically changing sets of heterogeneous processors and networks.

The Ibis distribution package comes with over 30 working applications, in the areas of physical simulations, parallel rendering, computational mathematics, state space search, bioinformatics, prime numbers factorization, data compression, cellular automata, grid methods, optimization, and generic problem solving. The Ibis applications closely resemble real-life parallel applications, as they cover a wide-range of computation/communication ratios, have different communication patterns and memory requirements, and are parameterized. Many of the Ibis

applications report detailed performance results. Last but not least, all the Ibis applications have been thoroughly described and tested [5, 16]. For a complete list of publications, please visit <http://www.cs.vu.nl/ibis>. Therefore, the Ibis applications are representative for grid applications written in Java, and can be easily included in synthetic grid workloads.

3.2 GRENCHMARK: Synthetic Grid Workloads

GRENCHMARK is a synthetic grid workload generator and submitter. It is *extensible*, in that it allows new types of grid applications to be included in the workload generation, *parameterizable*, as it allows the user to parameterize the workloads generation and submission, and *portable*, as its reference implementation is written in Python.

The workload generator is based on the concepts of *unit generators* and of job description files (JDF) *printers*. The *unit generators* produce detailed descriptions on running a set of applications (*workload unit*), according to the workload description provided by the user. In principle, there is one unit for each type of supported application type. The *printers* take the generated workload units and create job description files suitable for grid submission. In this way, multiple unit generators can be coupled to produce a workload that can be submitted to any grid resource manager, as long as the resource manager supports that type of applications.

The grid applications currently supported by GRENCHMARK are sequential jobs, jobs which use MPI, and Ibis jobs. We use the Ibis applications included in the default Ibis distribution (see Section 3.1). We have also implemented three *synthetic applications*: **sser**, a sequential application with parameterizable computation and memory requirements, **sserio**, a sequential application with parameterizable computation and I/O requirements, and **smi1**, an MPI application with parameterizable computation, communication, memory, and I/O requirements. Currently, GRENCHMARK can submit jobs to KOALA, Globus GRAM, and Condor.

The workload generation is also dependent on the applications inter-arrival time [7]. Peak job arrival rates for a grid system can also be modeled using well-known statistical distributions [7, 11]. Besides the Poisson distribution, used traditionally in queue-based systems simulation, modeling could rely on uniform, normal, exponential and hyper-exponential, Weibull, log normal, and gamma distributions. All these distributions are supported by the GRENCHMARK generator.

The workload submitter generates detailed reports of the submission process. The reports include all job submission commands, the turnaround time of each job, including the grid overhead, the total turnaround time of the workload, and various statistical information.

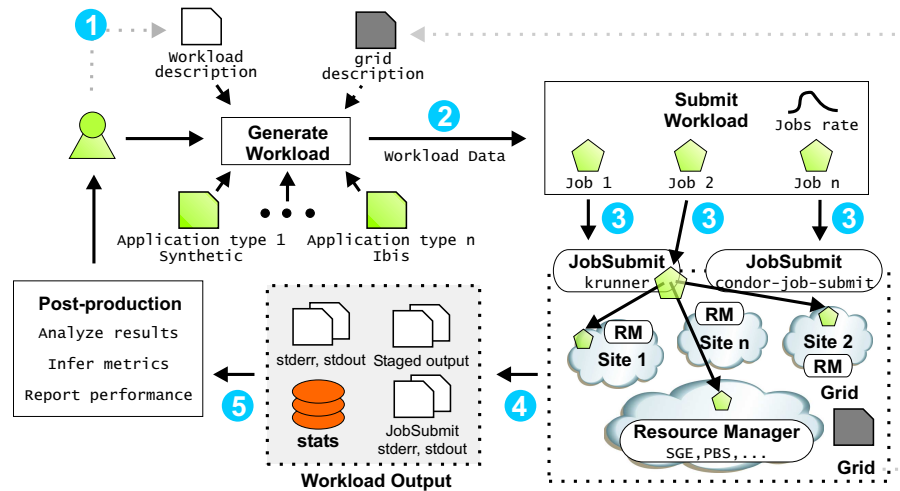


Fig. 1. The GRENCHMARK process.

3.3 Using the Framework

Figure 1 depicts the typical usage of our framework. First, the user describes the workload to be generated, as a formatted text file (1). Based on the user description, on the known application types, and on information about the grid sites, a workload is then generated by GRENCHMARK (2), and submitted to the grid (3). The grid environment is responsible for executing the jobs and returning their results (4). The results include not only job outcomes, but also detailed submission reports. Finally, the user processes all results in a post-production step (5).

4 A Concrete Case: Synthetic Workloads for the DAS

This section presents a concrete case for our framework: generating and running synthetic workloads on the DAS [1]. The Ibis applications were combined with the synthetic applications, to create a pool of over 35 grid applications. The GRENCHMARK tools were used to generate and launch the synthetic workloads.

4.1 KOALA: Scheduling Grid Applications

A key part of the experimental infrastructure is the KOALA [12] grid scheduler. To the author’s knowledge, KOALA is the only fault-tolerant, well-tested, and deployed grid scheduler that provides support for *co-allocated* jobs, that is, it can simultaneously allocate resources in multiple grid sites to single applications which consist of multiple components. KOALA was used to submit the generated workloads to the DAS multi-cluster. Its excellent reporting capabilities were also used for evaluating the jobs execution results.

Workload	Applications types	# of Jobs	# of CPUs	# of Components	Component Size	Success Rate
gmark1	synthetic, sequential	100	1	1	1	97%
gmark+	synthetic, seq. & MPI	100	1-128	1-15	1-32	81%
ibis1	N Queens, Ibis	100	2-16	1-8	2-16	56%
ibis+	various, Ibis	100	2-32	1-8	2-16	53%
wl+all	all types	100	1-32	1-8	1-32	90%

Table 1. The experimental workloads. As the DAS has only 5 sites; jobs with more than 5 components will have several components running at the same site.

```
# File-type: text/wl-spec
#ID Jobs Type SiteType Total SiteInfo ArrivalTimeDistr otherInfo
? 25 sser single 1 *:? Poisson(120s) StartAt=0s
? 25 sserio single 1 *:? Poisson(120s) StartAt=60s
? 25 smpi1 single 1 *:? Poisson(120s) StartAt=30s,ExternalFile=smpi1.xin
? 25 smpi1 single 1 *:? Poisson(120s) StartAt=90s,ExternalFile=smpi2.xin
```

Fig. 2. A GRENCHMARK workload description example.

For co-allocated jobs, KOALA gives the user the option to specify or not the actual execution sites, i.e., the clusters where job components should run. KOALA supports *fixed* jobs, for which users fully specify the execution sites, *non-fixed* jobs, for which the user does not specify the execution sites, leaving instead KOALA to select the best sites, and *semi-fixed* jobs, which are a mix of the previous two. KOALA may schedule different components of a non-fixed or of a semi-fixed job onto the same site. We used this feature heavily for the Ibis and synthetic MPI applications.

4.2 Workload Generation

Table 1 shows the structure of the five generated workloads, each comprising 100 jobs. To satisfy typical grid situations, jobs request resources from 1 to 15 sites. For parallel jobs, there is a preference for 2 and 4 sites. Site requests are either precise (specifying the full name of a grid site) or non-specified (leaving the scheduler to decide). For multi-site jobs, components occupy between 2 and 32 processors, with a preference for 2, 4, and 16 processors. We used combinations of parameters that would keep the run-time of the applications under 30 minutes, under optimal conditions. Each job requests resources for a time below 15 minutes. Various inter-arrival time distributions are used, but the submission time of the last job of any workload is kept under two hours.

Figure 2 shows the workload description for generating the **gmark+** test, comprising 100 jobs of four different types. The first two lines are comments. The next two lines are used to generate sequential jobs of types **sser** and **sserio**, with default parameters. The final two lines are used to generate MPI jobs of type **smpi1**, with parameters specified in external files **smpi1.xin** and **smpi2.xin**. All four job types assume an arrival process with Poisson distribution, with a average rate of 1 job every 120 seconds. The first job of each type starts at a time specified in the workload description with the help of the **StartAt** tag.

Job name	Job type	Turnaround [s]			Runtime [s]			Run	Run+Success
		Avg.	Min	Max	Avg.	Min	Max		
sser	sequential	129	16	926	44	1	588	100%	97%
smpi1	MPI	332	21	1078	110	1	332	80%	85%
N Queens	Ibis	99	15	1835	31	1	201	66%	85%

Table 2. A summary of time and run/success percentages for different job types.

4.3 The Workload Submission

GRENCHMARK was used to submit the workloads. Each workload was submitted in the normal DAS working environment, thus being influenced by the background load generated by other DAS users. Some jobs could not finish in the time for which they requested resources, and were stopped automatically by the KOALA scheduler. This situation corresponds to users under-estimating applications’ runtimes. Each workload ran between the submission start time and 20 minutes after the submission of the last job. Thus, some jobs did not run, as not enough free resources were available during the time between their submission and the end of the workload run. This situation is typical for real working environments, and being able to run and stop the workload according to the user specifications shows some of the capabilities of GRENCHMARK.

5 Experimental results

This section presents an overview of the experimental results, and shows that workloads generated with GRENCHMARK can cover in practice a wide-range of run characteristics.

5.1 Performance Results

Table 1 shows the success rate for all five workloads (column *Success Rate*). A successful job is a job that gets its resources, runs, finishes, and returns all results within the time allowed for the workload. The lower performance of Ibis jobs (workload **ibis+**) when compared to all the others, is caused by the fact that the system was very busy at the time of testing, making the resource allocation particularly difficult. This situation cannot be prevented in large-scale environments, and cannot be addressed without special resource reservation rights.

The turnaround time of an application can vary greatly (see Table 2), due to different parameter settings, or to varying system load. The variations in the application runtimes are due to different parameter settings.

As expected, the percentage of the applications that are actually run (Table 2, column *Run*) depends heavily on the job size and system load. The success rate of jobs that *did* run shows little variation (Table 2, column *Run+Success*). The ability of GRENCHMARK to report percentages such as these enables future work on comparing of the success rate of co-allocated jobs, versus single-site jobs.

5.2 Dealing With Errors

Using the combined GRENCHMARK and KOALA reports, it was easy to identify errors at various levels in the submission and execution environment: the user, the scheduler, the local and the remote resource, and the application environment levels. For a better description of the error levels, and for a discussion about the difficulty of trapping and understanding errors, we refer the reader to the work of Thain and Livny [14].

We were able to identify bottlenecks in the grid infrastructure, and in particular in KOALA, which was one of our goals. For example, we found that for large jobs in a busy system, the percentage of unsuccessful jobs increases dramatically. The reason is twofold. First, using a single machine to submit jobs (a typical grid usage scenario) incurs a high level of memory occupancy, especially with many jobs waiting for the needed resources. A possible solution is to allow a single KOALA job submitter to support multiple job submissions. Second, there are cases when jobs attempt to claim the resources allocated by the scheduler, but for some reason fail to do so. These jobs should not be re-scheduled immediately, or this could lead to a high occupancy of the system resources. A possible solution is to use an exponential back-off mechanism when scheduling such jobs.

6 Conclusions and Ongoing Work

This work has addressed the problem of synthetic grid workload generation and submission. We have integrated three research prototypes, namely a grid application development toolkit, Ibis, a grid metascheduler, KOALA, and a synthetic grid workload generator, GRENCHMARK, and used them to generate and run workloads comprising well-established and new grid applications on a multi-cluster grid. We have run a large number of application instances, and presented overview results of the runs.

We are currently adding to GRENCHMARK the complex applications generation capabilities and an automatic results analyzer. For the future, we plan to prove the applicability of GRENCHMARK for specific grid performance evaluation, such as such as an evaluation of the DAS support and performance for CERN's High-Energy Physics applications [6], or a performance comparison of co-allocated and single site applications, to complement our previous simulation work [3].

7 Acknowledgements

This research work is carried out under the FP6 Network of Excellence Core-GRID funded by the European Commission (Contract IST-2002-004265).

We would also like to thank Hashim Mohamed and Wouter Lammers, for their work on KOALA, and Gosia Wrzesiska, Niels Drost, and Mathijs den Burger, for their work on Ibis.

References

- [1] Henri E. Bal et al. The distributed ASCI supercomputer project. *Operating Systems Review*, 34(4):76–96, October 2000.
- [2] F. Berman, A. Hey, and G. Fox. *Grid Computing: Making The Global Infrastructure a Reality*. Wiley Publishing House, 2003. ISBN: 0-470-85319-0.
- [3] A. I. D. Bucur and D. H. J. Epema. Trace-based simulations of processor co-allocation policies in multiclusters. In *Proc. of the 12th IEEE HPDC*, pages 70–79, Washington, DC, USA, 2003. IEEE Computer Society.
- [4] Greg Chun, Holly Dail, Henri Casanova, and Allan Snavely. Benchmark probes for grid assessment. In *IPDPS*. IEEE Computer Society, 2004.
- [5] Alexandre Denis, Olivier Aumage, Rutger Hofman, Kees Verstoep, Thilo Kielmann, and Henri E. Bal. Wide-area communication for grids: An integrated solution to connectivity, performance and security problems. In *13th International Symposium on High-Performance Distributed Computing (HPDC-13)*, pages 97–106, Honolulu, Hawaii, USA, June 2004.
- [6] D. Barberis et al. Common use cases for a high-energy physics common application layer for analysis. Report LHC-SC2-20-2002, LHC Grid Computing Project, October 2003.
- [7] Steve J. Chapin et al. Benchmarks and standards for the evaluation of parallel job schedulers. In Dror G. Feitelson and Larry Rudolph, editors, *Job Scheduling Strategies for Parallel Processing*, pages 67–90. Springer-Verlag, 1999.
- [8] Michael Frumkin and Rob F. Van der Wijngaart. Nas grid benchmarks: A tool for grid space exploration. *Cluster Computing*, 5(3):247–255, 2002.
- [9] Vladimir Getov and Thilo Kielmann, editors. *Component Models and Systems for Grid Applications*, volume 1 of *CoreGRID seroes*. Springer Verlag, June 2004. Proceedings of the Workshop on Component Models and Systems for Grid Applications held June 26, 2004 in Saint Malo, France.
- [10] M. Humphrey et al. State and events for web services: A comparison of five WS-Resource Framework and WS-Notification implementations. In *4th IEEE International Symposium on High Performance Distributed Computing (HPDC-14)*, Research Triangle Park, NC, USA, July 2005.
- [11] Uri Lublin and Dror G. Feitelson. The workload on parallel supercomputers: Modeling the characteristics of rigid jobs. *J. Parallel & Distributed Comput.*, 63(11):1105–1122, Nov 2003.
- [12] H.H. Mohamed and D.H.J. Epema. Experiences with the koala co-allocating scheduler in multiclusters. In *Proc. of the 5th IEEE/ACM Int'l Symp. on Cluster Computing and the GRID (CCGrid2005)*, Cardiff, UK, May 2005.
- [13] Allan Snavely, Greg Chun, Henri Casanova, Rob F. Van der Wijngaart, and Michael A. Frumkin. Benchmarks for grid computing: a review of ongoing efforts and future directions. *SIGMETRICS Perform. Eval. Rev.*, 30(4):27–32, 2003.
- [14] Douglas Thain and Miron Livny. Error scope on a computational grid: Theory and practice. In *Proc. of the 11th IEEE HPDC*, page 199, Washington, DC, USA, 2002. IEEE Computer Society.
- [15] G. Tsouloupas and M. D. Dikaiakos. GridBench: A workbench for grid benchmarking. In P. M. A. Sloot, A. G. Hoekstra, T. Priol, A. Reinefeld, and M. Bubak, editors, *EGC*, volume 3470 of *LNCS*, pages 211–225. Springer, 2005.
- [16] Rob V. van Nieuwpoort, J. Maassen, G. Wrzesinska, R. Hofman, C. Jacobs, T. Kielmann, and H. E. Bal. Ibis: a flexible and efficient java-based grid programming environment. *Concurrency & Computation: Practice & Experience.*, 17(7-8):1079–1107, June-July 2005.