

# Grid Resource Ranking Using Low-Level Performance Measurements<sup>\*</sup>

George Tsouloupas and Marios D. Dikaiiakos

Dept. of Computer Science,  
University of Cyprus  
1678, Nicosia, Cyprus  
{georget,mdd}@cs.ucy.ac.cy

**Abstract.** This paper outlines a feasible approach to ranking Grid resources based on an easily obtainable application-specific performance model utilizing low-level performance metrics. First, Grid resources are characterized using low-level performance metrics; Then the performance of a given application is associated to the low-level performance measurements via a *Ranking Function*; Finally, the Ranking Function is used to rank all available resources on the Grid with respect to the specific application at hand. We show that this approach yields accurate results.

## 1 Introduction

Matching between resource requests and resource offerings is one of the key considerations in Grid computing infrastructures. Currently, the implementation of matching is based on the *matchmaking* approach introduced by the Condor project [6], adapted to multi-domain environments and Globus, and extended to cover aspects such as data access and work-flow computations, interactive Grid computing, and multi-platform interoperability. Matchmaking produces a ranked list of resources that are compatible to the submitted resource requests. Ranking decisions are based on published information regarding the number of CPU's of each resource, their nominal speed, the nominal size of main memory, the number of free CPU's, available bandwidth, etc. This information is retrieved from Grid information services such as the Monitoring and Discovery Service of Globus.

This approach works well in cases where the main consideration of end-users is to allocate sufficient numbers of idle CPU's in order to achieve a high job-submission throughput with opportunistic scheduling. In several scenarios, however, reliance to matchmaking is not sufficient; for instance, when end-users wish to “shop around” for Grid computing resources before deciding where to deploy a high-performance computing application, or when Virtual Organization (VO) operators want to audit the real availability and configuration status of their providers' computing resources [4]. In such cases, the information published by

---

<sup>\*</sup> This work was supported in part by the European Commission through projects EGEE (contract INFISO-RI-031688) and g-Eclipse (contract 034327).

resource providers and Grid monitoring systems does not provide sufficient detail and accuracy. Grid users need instead the capability to interactively administer benchmarks and tests, retrieve and analyze performance metrics, and select resources of choice according to their application requirements. To provide Grid users with such a *test-driving* functionality, we designed and implemented *GridBench*, a framework for evaluating the performance of Grid resources interactively. GridBench facilitates the definition of parameterized execution of various probes on the Grid, while at the same time allowing for archival, retrieval, and analysis of results [9,10]. GridBench comes with a suite of open-source micro-benchmarks and application kernels, which were chosen to test key aspects of computer performance, either in isolation or collectively (CPU, memory hierarchy, network, etc.) [11].

In this paper, we present *SiteRank*, a component that we developed on top of GridBench to support the user-driven ranking of computational Grid resources. SiteRank enables Grid users to easily construct and adapt ranking functions that:

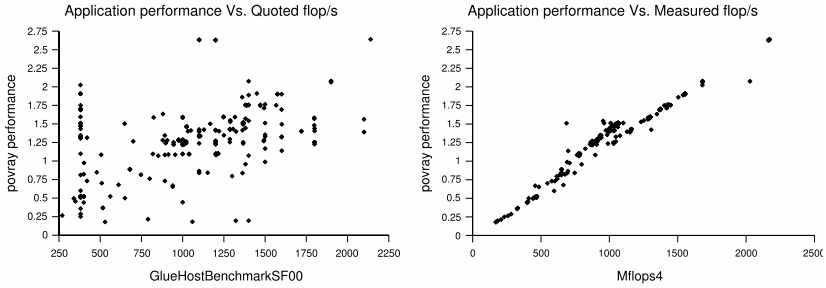
- (i) Take as arguments performance metrics derived with the low-level benchmarks of GridBench [11]; the selection of these metrics can be done manually or semi-automatically by the end-user, through the user interface of GridBench.
- (ii) Combine the selected metrics into a linear model that takes into account the particular requirements of the application that the user wishes to execute on the Grid (e.g., memory vs. floating-point performance bound). Using a ranking function, Grid users can derive rankings of Grid resources that are tailored to their specific application requirements.

In this paper, we describe the methodology followed by SiteRank to develop ranking functions. Furthermore, we demonstrate the use of SiteRank in the ranking of the computational resources of EGEE, which is the largest production-quality Grid in operation today [1]. To this end, we examine two alternative applications running on EGEE: povray, a ray-tracing application, and SimpleScalar, a simulator used for hardware-software co-verification and micro-architectural modelling. Our results show that SiteRank functions can provide an accurate ranking of EGEE resources, in accordance to the different requirements that each application has. Furthermore, that the careful selection of the low-level metrics used in the linear model is very important for the construction of accurate ranking functions.

The rest of this paper is organized as follows: Section 2 introduces SiteRank and its ranking methodology. Section 3 describes the use of SiteRank in the ranking of EGEE resources for the two applications of choice: povray and SimpleScalar. We conclude in Section 4.

## 2 SiteRank

Computational resources on the Grid exhibit considerable variance in terms of different performance characteristics. This leads to non-uniform application performance that significantly varies between applications.



**Fig. 1.** The relationship of application performance to **quoted** and **measured** metrics

One approach for ranking resources in terms of performance is the one taken by the current (EGEE) infrastructure, which is to publish GlueHostBenchmarkSF00 (SPEC-Float 2000 floating point performance metric) and GlueHostBenchmarkSI00 (SPEC-Int 2000 integer performance metric) values for each site. Unfortunately, values quoted by site administrators cannot be relied upon; This is evident in Figure 1 which compares the effectiveness of a **quoted** metric (Figure 1 left) in contrast to a measured metric (Figure 1 right). The charts speak for themselves; Clearly, the **quoted** metric does a very poor job in justifying application performance<sup>1</sup>. It is important to note that this would be *inadequate* even if the quoted values were correct, since application performance depends on much more than just two metrics (see Section 3).

## 2.1 The Ranking Methodology

The GridBench tool provides a *SiteRank module* that allows the user to interactively and semi-automatically build a *ranking model*. A *ranking model* consists of *filtering*, *aggregation* and *ranking functions* (Figure 2).

**Filtering** refers to a user selection regarding which results will be included or excluded in the ranking process. *Attribute filtering* allows the user to limit the selected set of measurements to the ones that match certain criteria in the benchmark description. E.g. limiting the selection to a specific VO, type of CPU, or the date and time results were obtained.

**Aggregation** allows the user to specify grouping of the measurements. The user can specify whether each measurement will count equally, irrespective of which worker-node it was executed on. In this case, the reported metric may possibly be less representative of the resource as a whole because some worker-nodes may be over-represented. On the other hand, this will tend to be more representative of what the user actually experiences once the resource's policy is applied. The *Aggregation* step produces a set of statistics for each metric: *mean*, *standard-deviation*, *min*, *max*, *average-deviation* and *count*.

<sup>1</sup> Similar results are obtained with GlueHostBenchmarkSI00 just as with GlueHostBenchmarkSF00.

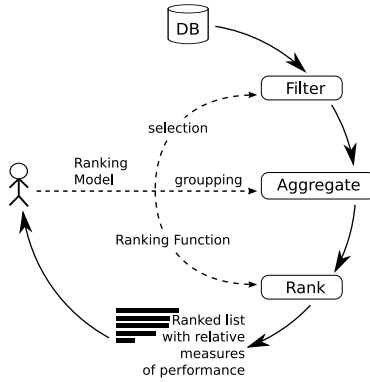
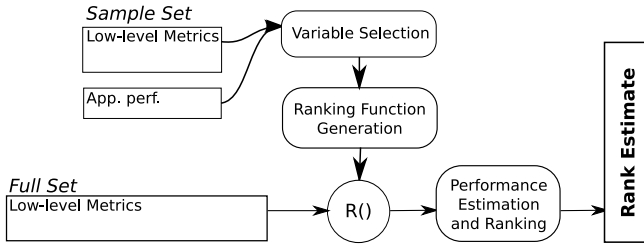


Fig. 2. The ranking process

During the aggregation step, the raw metrics are normalized according to a base value. The base values are configurable and in our experiments we used values from a typical 3.0GHz Xeon worker-node. For example, we used the value of 1050.0 to normalize the Mflops4 metric. The aggregation step is also important for the conversion of vector-type metrics, such as the ones produced by CacheBench into scalars (see later description on the *c512k* metric) so that they can be used in ranking functions.

**Ranking Function Construction:** The end goal of this methodology is a ranked list of computational resources that reflects the performance that users will experience running a specific application. It involves establishing a relationship between application performance and a set of low-level measurements. The process is illustrated in Figure 3, and it is outlined by the following steps (see Section 3 for an example):

1. **Sampling:** Obtain low-level performance metrics  $m$  for a small sample of resources – typically 10-15% of the full-set of resources. For the same sample of resources also obtain application performance measurements, i.e. application completion times. The application performance of this sample is denoted  $\alpha$  where each  $\alpha = 1/(\text{completion time})$ .
2. **Ranking Function Generation:** Determine a *Ranking Function*  $R$  based on the low-level metric data  $m$  and application performance  $\alpha$ , so that  $\alpha = R(m)$ . This involves the selection of the low-level metrics that closely correlate to this application’s performance, followed by a linear fit of the data, i.e. multivariate regression.
3. **Estimation:** For the set of the remaining resources, obtain only low-level performance metrics  $M$ , and apply the ranking function in order to obtain an estimate of the application performance  $A_{est}$  such that  $A_{est} = R(M)$ . Sorting  $A_{est}$  produces the *Rank Estimation*.



**Fig. 3.** Rank Estimate generation process outline

**Table 1.** Metrics and Benchmarks

Factor	Metric	Delivered By
CPU	Floating-Point operations per second	Flops
CPU	Integer operations per second	Dhrystone
Main memory	sustainable memory bandwidth in MB/s	Stream
Main memory	Available physical memory in MB	Memsize
Cache	memory bandwidth by varying array sizes in MB/s	CacheBench
Disk	Disk bandwidth for read/write/rewrite	bonnie++
Interconnect	latency, bandwidth and bisection bandwidth	MPPTest

## 2.2 Metrics

Selecting the *right* metrics to characterise the resources is of utmost importance in order to adequately characterize the major computational characteristics that affect application performance. In fact, we consider a *good set of metrics* one that can adequately explain the performance of several distinct applications. In the process of picking the right metrics and the right benchmarks to deliver these metrics, we limited ourselves to freely available tools that we could widely deploy and run. We also aimed at keeping the number of metrics low and we favored well-known metrics. A more detailed discussion can be found in [11].

Table 1 shows a list of low-level metrics and the associated benchmarks. The Flops benchmark yields 4 metrics, *Mflops1*, *Mflops2*, *Mflops3* and *Mflops4*, each consisting of different mixes of floating-point additions, subtractions multiplications and divisions. Dhrystone yields the *dhry* integer performance metric. The STREAM memory benchmark yields the *copy*, *add*, *multiply* and *triad* metrics which measure memory bandwidth using different operations. For cache metrics, measuring memory (cache) bandwidth  $B$  by allocating and accessing progressively larger array sizes  $s$ , the CacheBench benchmark produces a series of values  $B_s$  where  $s = 2^8, 2^9, 2^{10} \dots 2^n$ . By summing up the product of the bandwidths and respective sizes we derive a metric that takes into account both the cache size *and* the cache speed :  $\sum_{s=8}^n s \times B_s$ . For example, summing up to 512kb, i.e.  $\sum_{s=8}^{19} s \times B_s$  yields the *c512k* metric. This is done for sizes up to 512kb, 1Mb, 2Mb, 4Mb, 8Mb yielding the metrics *c512k*, *c1M*, *c2M*, *c4M* and *c8M* respectively. This approach alleviates the problem of looking up the cache size for the

multitude of CPU’s on the Grid, or detecting the cache sizes of a potentially multilevel cache.

### 3 Experimentation

In this section we demonstrate the proposed methodology by automatically determining a *Ranking Function*, obtaining a *Ranking Estimate* and validating that the Ranking Estimate is accurate by directly measuring the performance of the application. This is done for two applications, on a set of about 230 sites that belong to the EGEE infrastructure. We use two serial applications:

**povray:** The Povray v3.6 ray-tracing application using the a 40x40 scene.

**sisc:** The SimpleScalar, computer architecture simulation.<sup>2</sup>

For this experiment, we aimed at having between 2 and 3 measurements from each computational resource. One noteworthy fact is that we could only obtain results for about 160 out of the 230 sites. This was partly due to errors and site unavailability, but also due to exhausted quotas at some resources. We used the GridBench framework to obtain our measurements. The process of integrating the two applications into GridBench including the compilation took less than one hour and only needs to be performed once. The process of actually running all the experiments took less than 10 minutes, although we did have to wait for a few hours until the results from all the queued jobs were in. We then exported these results into an open-source statistics software package<sup>3</sup> (“R”).

The data-set obtained by running the benchmarks on all the available computational resources will be referred to from now on as the *full-set*. Out of the full-set, we obtained a random sample, henceforth referred to as the *sample-set*, with results from 24 resources (15% of the full-set). A *correlation matrix* indicates which metrics are most correlated to application performance; this is shown in Figure 4. The problem of collinearity must be taken into consideration when narrowing down the selection of metrics. As shown in Figure 4 some metric groups are highly collinear, in such cases we eliminate the collinear metrics by selecting one metric out of the group, i.e. the one with the highest correlation to the application. In this example we kept *Mflops4* and discarded *Mflops2*, *Mflops3* and *dhry*. Selecting the *Mflops4* and *c512k* metrics for building the Ranking Function, leads to the next step, i.e. calculating the *a* and *b* coefficients in order to best satisfy:

$$\alpha_{povray} = a \times Mflops4 + b \times c512k$$

Outlier removal is achieved by performing a linear regression, and data-points that fall more than two standard deviations away from the rest are filtered out. In our specific example, 2 out of the 18 points were dropped. Linear regression is

<sup>2</sup> Limited execution privileges for the Virtual Organization through which we performed our experiments, dictated that we use parameters resulting in short application completion times. This applied both to **povray** and to **sisc**.

<sup>3</sup> Use of the R software was limited to establishing the relationship between the low-level metrics and application performance, and the validation of the results. All charts included in the paper we created using GridBench.

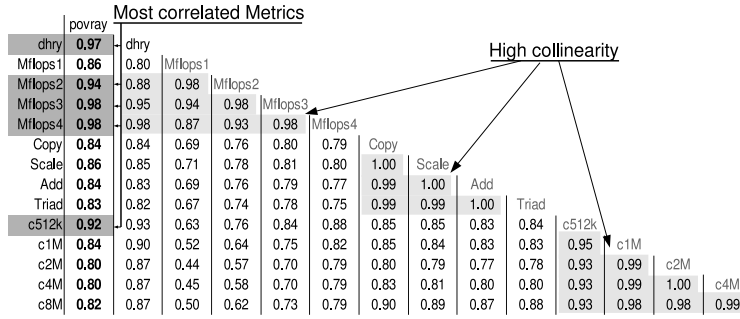


Fig. 4. Correlation Matrix for the *povray* application

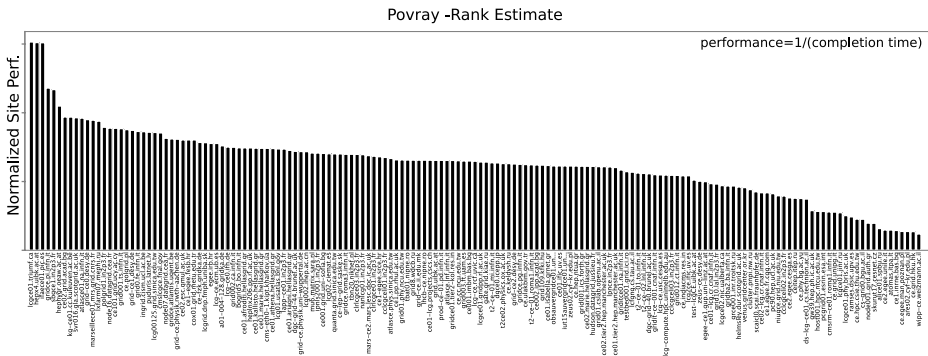


Fig. 5. Rank Estimate for the *povray* application

performed once again using the filtered sample-set, which yields the coefficients  $a = 0.94$  (for  $Mflops_4$ ) and  $b = 0.46$  (for  $c512k$ ). Finally, we apply this model on the *full-set* in order estimate the performance of the application:

$$A_{povray} = 0.94M_{Mflops_4} + 0.46M_{c512k}$$

Ordering the list of resources by  $A_{povray}$  gives the *Rank Estimate*. The Rank Estimate is shown in Figure 5. In order to test that the Ranking Estimate is accurate the performance of the application was directly measured for the whole infrastructure. This is only necessary in order to validate the model and not part of the methodology. The measured performance is shown in Figure 6. The agreement between the Rank Estimate and the measured ranking can be statistically tested by calculating the *rank correlation*. There are several ways of doing this, such as Kendall’s  $\tau$ , which ranges from -1 to 1 and is also known as the “bubble-sort distance”. Kendall’s  $\tau$  yielded  $\tau = 0.90$ . Spearman’s  $\rho$ , which again ranges from -1 to 1, yielded  $\rho = 0.977$ . Finally, Pearson’s correlation coefficient yielded 0.98. All three of the statistics show that the two rankings are quite similar. The  $\tau$  statistic appears considerably lower that the other two, due to the fact that our data-set contains a lot of resources that are of almost identical performance.

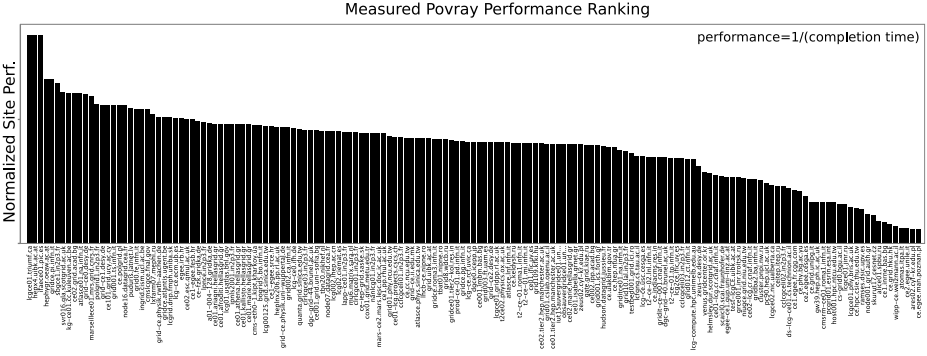


Fig. 6. Measured **povray** performance on 159 resources of the EGEE infrastructure

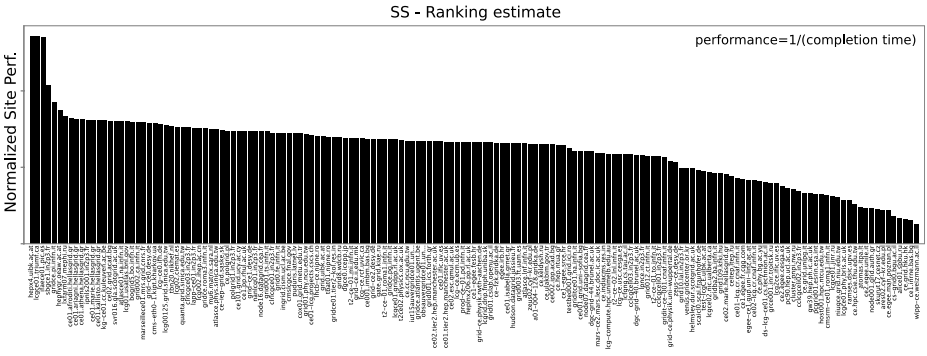


Fig. 7. Rank Estimate for the **sisc** application on the EGEE infrastructure

Extremely small fluctuations in measurement are enough to change the ordering. Yet, the performance of the resources is nearly identical, so the reordering is not very significant. For this reason the authors are inclined to take  $\rho = 0.977$  as the more representative measure.

For the second application **sisc** we used the same methodology and the same sample-set that was used in the previous case. The metrics dictated by the correlation matrix are *d hry* and *c512k*. Performing the regression, outlier removal and then estimating the metric coefficients yields:

$$A_{sisc} = 0.27M_{d hry} + 0.18M_{c512k}$$

The Ranking estimate is given in Figure 7. The correlation of estimated and actual is again quite high with a value of  $\rho = 0.959$ . Thus, for both applications the ranking of resources based on low-level measurements provides results that are very close to the ranking produced by running the application itself.

## 4 Conclusion

The work presented here, i.e. Ranking based on derived models of low-level metrics, describes an alternative way of choosing and ranking resources. We



propose a semi-automated user-driven approach to ranking Grid resources that employs user-specified metrics and *ranking functions*.

The process of running benchmarks collecting and analysing results and generating ranked lists, would simply not be feasible if it had to be done manually, especially if it had to be done by the end user. Eventually, resource performance information will be coupled with resource pricing information. Users will then be able to “shop around” and pick the right resources (e.g. black-listing or white-listing) in order to influence the matchmaking process in a way that benefits them. The SiteRank module of the GridBench tool allows the user to interactively construct and modify ranking functions based on the collected measurements. The *Ranking Estimate* has proven to be quite accurate with a very high correlation to measured application performance for at least two applications, *povray* and *SimpleScalar*.

We have illustrated that current approaches to expressing the performance of resources, such as publishing the *quoted*, not measured, GlueHostBenchmarkSF00 and GlueHostBenchmarkSI00 metrics into the information system are not satisfactory, since they do not correlate well with at least the two applications that we have investigated.

Other tools in the general area of Grid testing and benchmarking include the Grid Assessment Probes [3], DiPerF [5] and the Inca test harness and reporting framework [7]. These are testing/benchmarking frameworks that provide functionality ranging from testing of Grid services to the monitoring of service agreements. In contrast, we focus on user-driven performance exploration and ranking. Benchmarking as a data-source for resource-brokering is explored in [2]. This work suggests the application of *weights* to different resource attributes and the use of *application benchmarks* to obtain a ranking that can eventually be used for resource brokering; we have also suggested this in our previous work [8].

Choosing the right metrics to collect is of vital importance, as an incomplete set of metrics will yield poor characterization. For example, our initial experiments did not include metrics that characterize the memory cache. While we had been collecting measurements about the cache, the data was in a form that was rather difficult to integrate into a regular function. Also, we had falsely assumed that the cache effects would be largely accounted for in other metrics. The initial results were not at all encouraging; but including the cache metrics, i.e. *c512k*, completely changed the situation. Indicative was the improvement of the  $\rho$  rank correlation statistic from approximately  $\rho = 0.8$  to  $\rho = 0.96$  for the *SimpleScalar* application. This also confirms the importance of a well-sized, fast cache to computational applications.

Further plans include the investigation of more applications, especially applications that are not CPU/memory bound, in order to evaluate the extent to which the metrics that we collect provide sufficient characterization.

## References

1. Enabling Grids for E-Science project, <http://www.eu-egee.org/>
2. Afgan, E., Velusamy, V., Bangalore, P.V.: Grid resource broker using application benchmarking. In: Sloot, P.M.A., Hoekstra, A.G., Priol, T., Reinefeld, A., Bubak, M. (eds.) EGC 2005. LNCS, vol. 3470, pp. 691–701. Springer, Heidelberg (2005)

3. Chun, G., Dail, H., Casanova, H., Snavely, A.: Benchmark probes for grid assessment. In: 18th International Parallel and Distributed Processing Symposium (IPDPS 2004), CD-ROM / Abstracts Proceedings, 26-30 April 2004, Santa Fe, New Mexico, USA. IEEE Computer Society (2004)
4. Coles, J.: Grid Deployment and Operations: EGEE, LCG and GridPP. In: Proceedings of the UK e-Science All Hands Meeting 2005, (accessed October 2005) (2005), <http://www.allhands.org.uk/proceedings/2005>
5. Dumitrescu, C., Raicu, I., Ripeanu, M., Foster, I.: Diferf: an automated distributed performance testing framework. In: Proceedings of the 5th International Workshop on Grid Computing (GRID2004), IEEE Computer Society Press, Los Alamitos (2004)
6. Raman, R., Livny, M., Solomon, M.: Matchmaking: An extensible framework for distributed resource management. *Cluster Computing* 2(2), 129–138 (1999)
7. Smallen, S., Olschanowsky, C., Ericson, K., Beckman, P., Schopf, J.: The inca test harness and reporting framework. In: SC '04: Proceedings of the 2004. ACM/IEEE conference on Supercomputing, Washington, DC, USA, 2004, p. 55. IEEE Computer Society Press, Los Alamitos (2004)
8. Tiramo-Ramos, A., Tsouloupas, G., Dikaiakos, M.D., Sloot, P.: Grid Resource Selection by Application Benchmarking: a Computational Haemodynamics Case Study. In: Sunderam, V.S., van Albada, G.D., Sloot, P.M.A., Dongarra, J.J. (eds.) ICCS 2005. LNCS, vol. 3514, pp. 534–543. Springer, Heidelberg (2005)
9. Tsouloupas, G., Dikaiakos, M.D.: GridBench: A Tool for Benchmarking Grids. In: Proceedings of the 4th International Workshop on Grid Computing (Grid2003), pp. 60–67. IEEE Computer Society, Los Alamitos (2003)
10. Tsouloupas, G., Dikaiakos, M.D.: GridBench: A Workbench for Grid Benchmarking. In: Sloot, P.M.A., Hoekstra, A.G., Priol, T., Reinefeld, A., Bubak, M. (eds.) EGC 2005. LNCS, vol. 3470, pp. 211–225. Springer, Heidelberg (2005)
11. Tsouloupas, G., Dikaiakos, M.D.: Characterization of Computational Grid Resources Using Low-level Benchmarks. In: Second IEEE International Conference on e-Science and Grid Computing (e-Science'06), pp. 70–77. IEEE Computer Society, Los Alamitos (2006)