

## INFRASTRUCTURE-AGNOSTIC PROGRAMMING AND INTEROPERABLE EXECUTION IN HETEROGENEOUS GRIDS

Enric TEJEDOR

*Conseil Européen pour la Recherche Nucleaire (CERN)  
Meyrin, Switzerland  
e-mail: enric.tejedor.saavedra@cern.ch*

Javier ÁLVAREZ

*The University of Adelaide  
Adelaide, Australia  
e-mail: javier.alvarez@adelaide.edu.au*

Rosa M. BADIA

*Barcelona Supercomputing Center (BSC-CNS)  
Jordi Girona 29, 08034 Barcelona (Spain)  
&  
Artificial Intelligence Research Institute (IIIA)  
Spanish Council for Scientific Research (CSIC)  
E-08193 Bellaterra, Barcelona (Spain)  
e-mail: rosa.m.badia@bsc.es*

**Abstract.** In distributed environments, no matter the type of infrastructure (cluster, grid, cloud), portability of applications and interoperability are always a major concern. Such infrastructures have a high variety of characteristics, which brings a need for systems that abstract the application from the particular details of each infrastructure. In addition, managing parallelisation and distribution also complicates the work of the programmer. In that sense, this paper demonstrates how

an e-Science application can be easily developed with the COMPSs programming model and then parallelised in heterogeneous grids with the COMPSs runtime. With COMPSs, programs are developed in a totally-sequential way, while the user is only responsible for specifying their tasks, i.e. computations to be spawned asynchronously to the available resources. The COMPSs runtime deals with parallelisation and infrastructure management, so that the application is portable and agnostic of the underlying infrastructure.

**Keywords:** Grid programming models, workflow managers, parallelism exploitation

**Mathematics Subject Classification 2010:** 68-N19, 68-M14

## 1 INTRODUCTION

In distributed environments, no matter the type of infrastructure (cluster, grid, cloud), portability of applications and interoperability are always a major concern [3, 2]. Different infrastructures can have very diverse characteristics. Besides, even in the scope of a given infrastructure, there is typically a plethora of alternatives to implement and execute an application, and often several vendors compete to dominate the market. Choosing one of the alternatives usually ties the application to it, e.g. due to the use of a certain API. As a result, it may be hard to port the application, not only to another kind of infrastructure, but also to an equivalent platform provided by another vendor or managed by different software.

Standards do appear, either “de facto” or produced by collaborative organisations that develop them, as in the case of the Open Grid Forum [7], but it is often complicated for them to be widely accepted. This situation, which is likely to keep happening in future scenarios, increases the importance of systems that free the user from porting the same application over different platforms.

On the other hand, some of the difficulties of programming applications for distributed infrastructures are not related to their particular characteristics, but to the duty of parallelisation and distribution itself [22]. This includes aspects like thread creation and synchronisation, messaging, data partitioning and transfer, etc. Having to deal with such aspects can significantly complicate the work of the programmer as well.

In that sense, this paper demonstrates how the COMPSs programming model and runtime system can be used to easily develop and parallelise applications in distributed infrastructures. More precisely, we discuss an example of an e-Science application that was programmed with the COMPSs model. Such application does not include any API call, deployment or resource management detail that could tie it to a certain platform. In addition, the application is programmed in a fully-sequential

fashion, freeing the programmer from having to explicitly manage parallelisation and distribution.

Furthermore, we present some experiments that execute that application in large-scale heterogeneous grids controlled by different types of middleware. A runtime is responsible for hiding that heterogeneity to the programmer, interacting with the grids and making them interoperable to execute the application. Consequently, the application remains agnostic of the underlying infrastructure, which favours portability.

The paper is structured as follows. Section 2 provides an overview of the COMPSs programming model and runtime system. Section 3 introduces the use-case e-Science application. Section 4 describes the grid testbed used in the experiments. Section 5 presents the results of the experiments. Finally, Section 6 discusses some related work and Section 7 concludes the paper.

## 2 OVERVIEW OF COMP SUPERSCALAR

This section introduces the COMP Superscalar (COMPSs) programming model, as well as the runtime system that supports the model's features. COMPSs is tailored for Java applications running on distributed platforms like clusters, grids and clouds. For a more detailed description of COMPSs, please see [28, 30, 29].

### 2.1 Programming Model

The COMPSs programming model can be defined as task-based and dependency-aware. In COMPSs, the programmer is only required to select a set of methods and/or services called from a sequential Java application, for them to be run as *tasks* – asynchronous computations – on the available distributed resources.

The task selection is done by providing a Task Selection Interface (TSI), a Java interface which declares those methods/services, along with some metadata. Part of these metadata specifies the direction (input, output or in-out) of each task parameter; this is used to discover, at execution time, the data dependencies between tasks. The TSI is not a part of the application: it is completely separated from the application code and it is not implemented by any of the user's classes; its purpose is merely specifying the tasks.

With COMPSs, sequential Java applications can be parallelised with no modifications: the application code does not contain any parallel construct, API call or pragma. All the information needed for parallelization is contained in the TSI. Besides, the application is not tied to a particular infrastructure: it does not include any resource management or deployment information.

### 2.2 Runtime System

The runtime system receives as input the class files corresponding to the sequential code of the application and the TSI. Before executing the application, the run-

time transforms it into a modified bytecode that can be parallelised. In particular, the invocations of the user-selected methods/services are automatically replaced by an invocation to the runtime: such invocation will create an asynchronous task and let the main program continue its execution right away.

The created tasks are processed by the runtime, which dynamically discovers the dependencies between them, building a task dependency graph. Moreover, a renaming technique is used to avoid some kinds of dependencies. The parallelism exhibited by the graph is exploited as much as possible, scheduling the dependency-free tasks on the available resources. The scheduling is locality-aware: nodes can cache task data for later use, and a node that already has some or all the input data for a task gets more chances to run it.

The interaction of the runtime with the infrastructure is done through JavaGAT [16], which offers a uniform API to access different kinds of grid middleware. COMPSs uses JavaGAT for two main purposes: submitting tasks and transferring files to grid resources. Thus, the runtime is responsible for transferring task data and managing task execution through JavaGAT, while the application is totally unaware of such details.

### 3 THE SIMDYNAMICS APPLICATION

DISCRETE [12] is a package devised to simulate the dynamics of proteins using the Discrete Molecular Dynamics (DMD) methods. In such simulations, the particles are assumed to move with constant velocity until a collision occurs, conserving the total momentum and energy, which drastically saves computation time compared to standard MD protocols.

The simulation program of DISCRETE receives as input a coordinate and a topology files, which are generated with a setup program also included in the package. The coordinate file provides the position of each atom in the structure, and the topology file contains information about the chemical structure of the molecule and the charge of the atoms. Besides, the simulation program reads a parameter file, which basically specifies three values: EPS (Coulomb interactions), FSOLV (solvation) and FVDW (Van Der Waals terms).

The SimDynamics application, which will be used in the experiments presented in Section 5, is a sequential Java program that makes use of the DISCRETE package. Starting from a set of protein structures, the objective of SimDynamics is to find the values of the EPS, FSOLV and FVDW parameters that minimise the overall energy obtained when simulating their molecular dynamics with DISCRETE. Hence, SimDynamics is an example of a parameter-sweeping application: for each parameter, a fixed number of values within a range is considered and a set of simulations (one per structure) is performed for each combination of these values (configuration). Once all the simulations for a specific configuration have completed, the configuration's score is calculated and later compared to the others in order to find the best one.

The main program of the SimDynamics application is divided in three phases:

1. For each of the  $N$  input protein structures, their corresponding topology and coordinate files are generated. These files are independent of the values of EPS, FSOLV and FVDW.
2. Parameter-sweep simulations: a simulation is executed for each configuration and each structure. These simulations do not depend on each other. The more values evaluated for each parameter, the more accurate will be the solution.
3. Finding the configuration with minimal energy: the execution of each simulation outputs a trajectory and an energy file, which are used to calculate a coefficient for each configuration. The main result of the application is the configuration that minimises that coefficient.

In order to run SimDynamics with COMPSs, a total of six methods invoked from the application were chosen as tasks. This was done by defining a TSI that declares those methods. Figure 1 contains a fragment of this TSI, more precisely the selection of method `simulate` as a task. The parameters of `simulate` are three input files, an input string and an output file. The declarations of the other five methods are analogous to this one. The following points describe the method tasks, the subindexes indicating the phase to which they belong:

- *genReceptorLigand*<sub>1</sub>: given a structure file, it generates some associated files (receptor and ligand). It is invoked  $N$  times (one per structure).
- *dmdSetup*<sub>1</sub>: it executes the DMDSetup binary, included in the DISCRETE package, with a structure's receptor and ligand as input; as output, it generates the topology and coordinate files for the structure. It is invoked  $N$  times (one per structure).
- *simulate*<sub>2</sub>: it runs the simulation binary of the DISCRETE suite, given a coordinate file, a topology and a specific configuration (FVDW, FSOLV and EPS values); it returns an average score file. If the number of values considered for EPS, FSOLV and FVDW is  $S_{EPS}$ ,  $S_{FSOLV}$  and  $S_{FVDW}$ , respectively, this method is invoked  $N \times S_{EPS} \times S_{FSOLV} \times S_{FVDW}$  times.
- *merge*<sub>2</sub>: it merges two average score files belonging to the same configuration of parameters. It is invoked  $(N - 1) \times S_{EPS} \times S_{FSOLV} \times S_{FVDW}$  times.
- *evaluate*<sub>3</sub>: it generates the final coefficient from all the average scores of a configuration. It is invoked once per configuration, i.e.  $S_{EPS} \times S_{FSOLV} \times S_{FVDW}$  times.
- *min*<sub>3</sub>: it receives two coefficient files and outputs the lowest one. It is invoked  $(S_{EPS} \times S_{FSOLV} \times S_{FVDW}) - 1$  times.

```

public interface SimDynamicsItf {
    @Method(declaringClass = "simdynamics.SimDynamicsImpl")
    void simulate(
        @Parameter(type = FILE) String paramFile,
        @Parameter(type = FILE) String topFile,
        @Parameter(type = FILE) String crdFile,
        String natom,
        @Parameter(type = FILE, direction = OUT) String average
    );
    ...
}

```

Figure 1. Code snippet of the Task Selection Interface for the SimDynamics application, where the `simulate` method is selected as a task. The `@Method` annotation specifies the class that implements `simulate`, and the `@Parameter` annotation contains parameter-related metadata (type, direction).

## 4 TESTBED INFRASTRUCTURE

The SimDynamics application was executed with COMPSs on real large-scale scientific grids. The whole infrastructure used in the tests is depicted in Figure 2, and it includes three grids: the Open Science Grid, Ibergrid and a small grid owned by the Barcelona Supercomputing Center [1].

Such infrastructure represents a heterogeneous testbed, comprised by three grids belonging to different administrative domains and managed by different middleware. The next subsections briefly describe the topology of these grids and explain how the COMPSs runtime was able to hide the complexity of their heterogeneity, keeping the grid-related details transparent to the application.

### 4.1 Grids

#### 4.1.1 Open Science Grid

The Open Science Grid (OSG) [9] is a science consortium, funded by the United States Department of Energy and the National Science Foundation, that offers an open grid cyberinfrastructure to the research and academic communities. OSG federates more than 100 sites around the world, most of them located in the United States, including laboratory, campus, and community facilities. These sites provide guaranteed and opportunistic access to shared computing and storage resources. As of May 2011, OSG comprised a total of around 70 000 cores and 29 Petabytes of disk storage and it provided 1.4 million CPU hours/day [10].

OSG is used by scientists and researchers to perform data analysis tasks that are too computationally intensive for a single data center or supercomputer. This grid was created to process data coming from the Large Hadron Collider at CERN, and consequently most of its resources are allocated for particle physics; however,

it is also used by research teams from disciplines like biology, chemistry, astronomy and geographic information systems.

Each of the OSG sites – clusters, computing farms – is configured to deploy a set of grid services, like user authorisation, job submission and storage management. Basically, a site is organised in a *Compute Element* (CE), running in a front-end node known as the *gatekeeper*, plus several *worker nodes* (or execution nodes). The CE allows users to run jobs on a site by means of the Globus GRAM (Grid Resource Allocation Manager) [19] interface; at the back-end of this GRAM gatekeeper, each site features one or more local batch systems – like Condor [31], PBS [8] or LSF [11] – that process a queue of jobs and schedule them on the worker nodes. Besides, the standard CE installation includes a GridFTP server; typically, the files uploaded to this server are accessible from all the nodes of the site via a distributed file system like NFS (Network File System [6]).

#### 4.1.2 Ibergrid

Ibergrid was set up in May 2010 as an umbrella organisation for ES-NGI [13] and INGRID [5] – the Spanish and Portuguese National Grid Initiatives, respectively – in the framework of the European Grid Initiative (EGI), which has the mission of creating and maintaining a pan-European Grid infrastructure.

Ibergrid offers aggregated computing power of more than 24,000 cores and 20 Petabytes of online storage and supports scientists in several fields of research, including high-energy physics, computational chemistry, engineering and nuclear fusion. Ibergrid also dedicates, like the OSG, a significant part of its resources to process data from the LHC. In total, the usage of Ibergrid reached 124 million CPU hours in 2011 [4].

Similarly to OSG, the Ibergrid infrastructure is composed by different sites, each one with a gatekeeper node interfacing to the cluster, a local resource management system (batch) and a set of worker nodes. However, in Ibergrid the middleware installed is gLite [21] and job management is a bit different: instead of submitting the jobs to a given CE directly, the user proceeds by interacting with a *Workload Management Server* (WMS), which acts as a meta-scheduling server. Therefore, matchmaking is performed at a higher level: the WMS interrogates the Information Supermarket (an internal cache of information) to determine the status of computational and storage resources, and the File Catalogue to find the location of any required input files; based on that information, the WMS selects a CE where to execute the job.

#### 4.1.3 BSC Grid

Finally, the Barcelona Supercomputing Center (BSC) Grid [1] is a small cluster located in the BSC premises and formed by five nodes. Three of them have a single-core processor at 3.60 GHz, 1 GB of RAM and a local disk of 60 GB. The other two have a quad-core processor at 2.50 GHz each core, 4 GB of RAM and

a local disk of 260 GB. The cluster does not have any shared file system configured.

The BSC Grid is the only grid of the testbed that supports interactive execution: the user can connect to any of the nodes separately via SSH and launch computations on them. Moreover, files can be transferred to/from the local disk of each node through SSH as well.

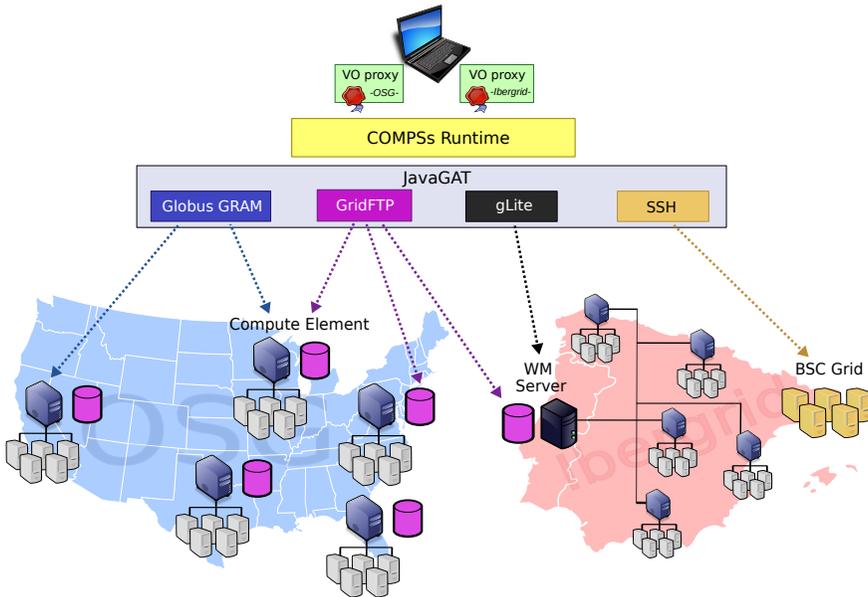


Figure 2. Testbed comprising two large-scale scientific grids (Open Science Grid, Ibergrid) and a BSC-owned grid. The SimDynamics application, running on a machine with COMPSs, interacts with the grids through JavaGAT and its middleware adaptors.

## 4.2 Configuration and Operation Details

In order to run the SimDynamics application in the described testbed, the testing environment was configured as shown in Figure 2.

The access point to the Grid was a laptop equipped with a dual-core 2.8 GHz processor and 8 GB RAM. This is different from the traditional procedure of submitting jobs from a User Interface node (UI) of a grid, where the software to interact with that grid is already present. Since the experiments did not target a particular grid but three different ones, and to illustrate how a user can execute a COMPSs application on the grid from his/her own machine, another approach was followed.

The laptop hosted the main program of the application, and therefore it had the COMPSs runtime and the JavaGAT library installed. Notice that no client

middleware had to be installed in the laptop, the GAT adaptors sufficed to interact with all the grids. In addition, prior to the execution, the credentials for each grid were obtained. Putting aside the setup of the SSH keys to access the BSC Grid, OSG and Ibergrid required proxy certificates for authentication, each with a different VO extension. Both proxies were created in a UI node of Ibergrid with the VOMS tools [14] and then placed in the laptop, so that JavaGAT could make use of those credentials when contacting the grids.

Concerning the grid middleware, the points below list the GAT adaptors and the corresponding grids where they were used. The resources available in each grid were specified in a resources file, along with their capabilities (e.g. associated storage servers).

- *Globus GRAM and OSG*: a total of six OSG sites that support our VO (Engage) were used in the tests, each with its own CE. The gatekeeper of every CE was contacted by means of the Globus GRAM adaptor, used for job submission and monitoring in OSG.
- *gLite and Ibergrid*: the gLite adaptor was used to submit and monitor jobs by connecting to an Ibergrid WMS, which is in charge of selecting the execution site in Ibergrid. Among all the WMS at the disposal of our VO (ICT), the one with most availability was chosen.
- *GridFTP (OSG and Ibergrid)*: the OSG CEs and the Ibergrid WMS offer each a GridFTP server. The GAT GridFTP adaptor was used to transfer files to those servers during execution.
- *SSH and BSC Grid*: two nodes of BSC Grid were used in the tests, being accessed through the GAT SSH adaptors, more precisely for job submission and file transfer.

Before execution, there was a previous phase of deployment where some required files were installed in the grids; those included, on the one hand, the worker runtime and, on the other, the classes and executables of the application tasks. In OSG, the files to be deployed were copied to the GridFTP server of each CE, so they could be accessed from the worker nodes. In Ibergrid, the files were transferred to the GridFTP server of the WMS, since the final execution site is not known in advance in this scenario; each time a job is created in Ibergrid, those files are copied by the worker runtime from the GridFTP server to the site where the job will run. Finally, in BSC Grid the files were placed in the local disk of the nodes.

At execution time, the master runtime of COMPSs sends the SimDynamics tasks and transfers files to the three grids by means of GAT. In OSG, the input files of each task are first pre-staged to the GridFTP server of the target CE, thus being accessible through the NFS server of that CE too; after that, when the job is created in the CE to execute the task, the worker runtime copies the input files from NFS to the local disk of the target worker node; similarly, the output files are copied from local to NFS at the end of the task, thus being available in the GridFTP server as well. In Ibergrid, the task input files are transferred to the GridFTP server of the

WMS; the pre- and post-staging of those files to/from the final worker node is taken care by gLite: the WMS chooses the execution site, sends the job to the head node of that site, then the job is locally scheduled and the input files are copied from the GridFTP server to the local disk of the worker node (the process is inverse for the output files). Lastly, the BSC Grid scenario is simpler since the files can be directly transferred to/from the local disk of the final execution node.

In the case of SimDynamics, all the application input files were initially located in the laptop's disk and then progressively transferred to the execution resources as the application ran; nevertheless, for applications dealing with huge files, the programmer can also refer to those files with a whole URI (i.e. including the resource name) in the application code, so that they are gotten from that resource.

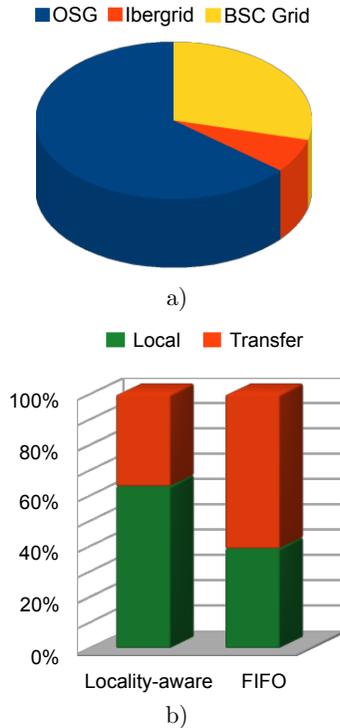
When scheduling jobs on the grids, the COMPSs runtime takes into account locality: a task will be assigned, if possible, to a resource that already possesses one or more of the task's input files (in its GridFTP server or local disk). Whenever a resource is freed (a task finishes), the scheduler chooses the task with the best score among the pending ones, the score being the number of task input files in the resource. Note that Ibergrid counts as a single entity for locality, because the final destination of the job is not decided by COMPSs. If some input file is missing in the chosen resource, such file is replicated to that resource. If the source and destination resources share the same credentials (e.g. two OSG sites) such transfer happens directly between them; otherwise, the file is first copied to the laptop and then to the destination resource.

## 5 EVALUATION

This section presents the results of executing the SimDynamics application (Section 3) in the described testbed (Section 4). These tests will show how the tasks of an e-Science application are executed in three different grids with COMPSs.

From the point of view of the application, *all the Grid management* discussed in Section 4 *is transparent*. The application deals with its parameters, i.e. number of structures and coefficients. For these experiments, the parameters were the following:  $N = 10$  (structures),  $S_{EPS} = 3$ ,  $S_{FSOLV} = 3$ ,  $S_{FVDW} = 3$  (i.e. 27 configurations for parameter sweeping). Applying the formulae in Section 3, this leads to a total of 586 tasks, including 270 simulation tasks – the most computationally-intensive with about two minutes of execution time each. The rest of the tasks are lightweight, with a duration of less than 10 seconds.

Figure 3 a) shows how tasks were distributed among the three grids during an execution of SimDynamics with COMPSs. The six OSG resources were the ones that consumed more tasks; indeed, among all the OSG sites that support our VO, the ones with most availability were chosen. The two BSC Grid nodes also executed a considerable number of tasks because they are directly accessible and therefore those tasks did not suffer from queue waiting times. Ibergrid received less load because of three factors. First, the Ibergrid queue times in these tests were high, which



caused tasks scheduled in Ibergrid to wait. Second, regarding the internal scheduling policies of the Ibergrid sites, several sites offer to our VO only opportunistic access to their resources; some other sites reserve a certain number of slots with priority but they are shared by all the Ibergrid VOs. Finally, the errors when submitting tasks to the WMS were quite frequent, which made tasks go through a (sometimes long) resubmission process.

In that sense, Table 1 contains the statistics of errors in task submissions and file transfers for the different grids and a particularly faulty execution of SimDynamics, in order to demonstrate the fault tolerance mechanisms of the COMPSs runtime. In general, the OSG sites presented only occasional failures in task submissions and file transfers, which were easily solved with resubmissions and retransfers with no need for task rescheduling. On the contrary, the errors when connecting to the Ibergrid WMS were common, possibly because of a bug in the JavaGAT gLite adaptor or because of the WMS itself; in order to face that issue, several retries were attempted when necessary for a task (6 per task on average), progressively increasing the time between two resubmissions. The most reliable combination of grid/adaptor was BSC Grid/SSH, for which no errors of any kind were registered.

Regarding data locality, Figure 3 b) illustrates the benefits of using a locality-aware task scheduling algorithm. Such algorithm is especially important in a highly-

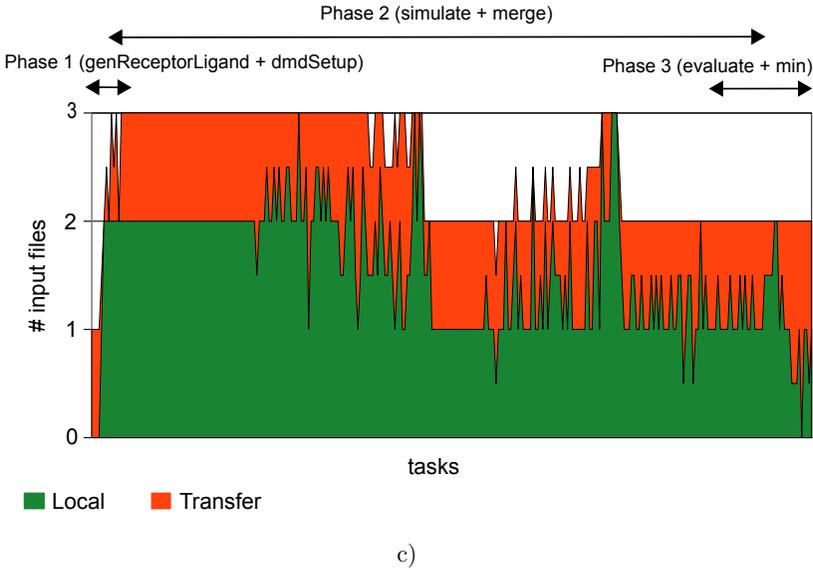


Figure 3. Test results for the SimDynamics application when run with COMPSs in the Grid testbed: a) distribution of the SimDynamics tasks among the three grids; b) comparison of percentage of transfers between the locality-aware and FIFO scheduling algorithms; c) evolution of the number of transfers when applying locality-aware scheduling.

Grid	Resource	# Task sub.		# File tra.	
		OK	Failed	OK	Failed
OSG	brgw1.renci.org	72	4	102	1
	gridgk01.racf.bnl.gov	43	0	70	1
	rossmann-osg.rcac.purdue.edu	57	14	89	11
	smufarm.physics.smu.edu	69	1	92	1
	stargrid02.rcf.bnl.gov	55	0	90	1
	u2-grid.ccr.buffalo.edu	62	1	96	0
	<i>TOTAL</i>	358	20	539	15
Ibergrid	wms01.ific.uv.es	33	209	58	0
	<i>TOTAL</i>	33	209	58	0
BSC Grid	bscgrid05.bsc.es	122	0	116	0
	bscgrid06.bsc.es	73	0	79	0
	<i>TOTAL</i>	195	0	195	0
	<b>TOTAL</b>	586	229	792	15

Table 1. Task submission and file transfer statistics for SimDynamics

distributed testbed like the one in Figure 2, where data transfers are costly. Figure 3 b) compares two executions of SimDynamics, one using locality-aware scheduling and another one applying a FIFO (First In First Out) strategy, and it shows the percentage of transfers actually performed versus the percentage of locality (the transfer was not necessary because the input file was already on the target execution resource), the total being the number of input files of all tasks. The locality-aware algorithm achieved remarkable results, preventing almost 2 out of every 3 transfers.

As a complement to Figure 3 b), Figure 3 c) illustrates the number of transfers that could be avoided thanks to locality all along the application execution. The  $x$ -axis represents the tasks of SimDynamics in the order that they are scheduled during the application execution; each point of the axis corresponds to two tasks, so that the number of points is reduced by half and the shape of the plotted lines is smoother. The  $y$ -axis reflects the evolution of avoided transfers (Local) and performed transfers (Transfer), each point showing the average of two tasks for both values. In the first phase of the application, the `genReceptorLigand` tasks require their input files to be transferred from the laptop to the grid resources, while the successor tasks `dmdSetup` benefit from full locality because they are scheduled in the same resources as their predecessor tasks, where the corresponding receptor and ligand files are already present. After that, there is an explosion of, first, `simulate` and, later, `merge` tasks, for which the runtime can prevent up to three and two transfers, respectively. Finally, the graph gets narrower when the merged scores of the simulations are processed by the `evaluate` and `min` tasks, each with two input files subject to locality.

```

public interface SimDynamicsIltf {
    @Constraints(operatingSystem = "Scientific Linux")
    @Method(...)
    void genReceptorLigand(...);

    @Constraints(appSoftware = "DISCRETE")
    @Method(...)
    void simulate(...);

    @Constraints(memory = 4)
    @Method(...)
    void evaluate(...);

    ...
}

```

Figure 4. Detail of the task constraint specification in the TSI of SimDynamics

SimDynamics works with only a few MB of data, but preventing files from being transferred in grids becomes more important as the size of these data increases. Furthermore, when dealing with big files the locality algorithm should take into account not only the number of files but also their size when selecting the destination host of a task. This requires to keep track of the sizes of each file updated/generated in the workers, as well as to send that information to the master runtime for it to

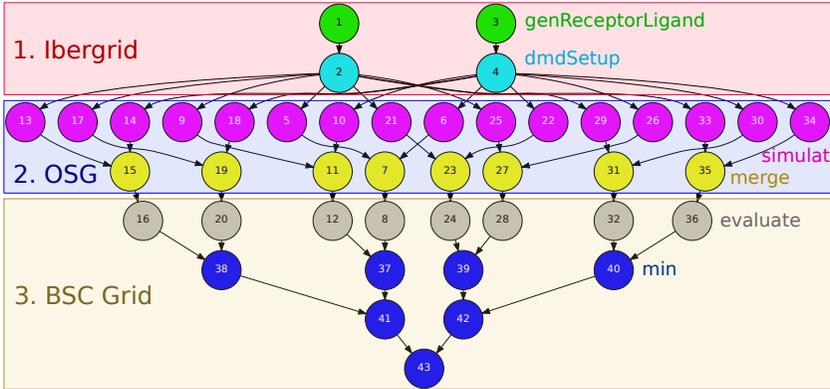


Figure 5. Reduced version of the SimDynamics graph (the real one contains 586 tasks). The constraints in Figure 4 lead to the task scheduling on the grids represented by this figure.

make better decisions. Such optimisation was addressed in [25] but it is out of the scope of this paper. Alternatively, the user can associate a given kind of task that accesses some big input data with a certain resource that is known to host those data, or with a resource that fulfills some other hardware/software requirements of the task.

In that sense, a variant of the tests discussed above intended to demonstrate how to use constraints to force the scheduling of tasks on certain resources. Let us assume that each kind of task in SimDynamics has some resource requirements; Figure 4 shows how they can be specified in COMPSs by means of the @Constraints annotation, at method level, in the TSI. Those requirements need to match the resource capabilities contained in the resources file. In this example, `genReceptorLigand` must be executed in nodes running Scientific Linux, which is the operating system installed in Ibergrid. Second, `simulate` is supposed to run in resources where the DISCRETE software is present; here, such capability was assigned only to OSG sites. Finally, `evaluate` has a hardware constraint attached – more precisely, the amount of physical memory – which was only known and specified in the resources file for the BSC Grid nodes. The three other kinds of task not shown in Figure 4 have analogous constraints.

As a result of the constraints, at execution time the scheduling of tasks on resources was the one depicted in Figure 5. This graph is a smaller version (only 8 configurations) just for illustration purposes. In conclusion, the programmer can use task constraints to make sure that a given group of tasks will be executed in one or more resources that conform to a set of requirements.

## 6 RELATED WORK

Apart from COMPSs, there exist other programming models for Grid applications. Ninf-G [26] offers a programming model where client programs can call libraries on remote resources using a client API that is built on top of the Globus Toolkit. Ninf-G's model is more complex than COMPSs', since the programmer has to substantially modify the original application code by including the invocations to the GridRPC API. Furthermore, COMPSs can submit tasks using different kinds of grid middleware. Satin [32] permits to express divide-and-conquer parallelism in Java applications, marking method invocations for asynchronous spawning. Nevertheless, the programmer must explicitly use a synchronisation primitive to wait for the spawned tasks; unlike Satin, COMPSs takes care of task and data synchronisation automatically, and it is not restricted to the divide-and-conquer paradigm. OpenWP [17] is a grid programming and runtime environment with a set of directives that have to be included in the application code to express parallelism and distribution. The main difference between COMPSs and OpenWP is that the latter requires to indicate the dependencies between tasks in the application code, whereas the former finds them automatically at execution time. ASSIST [15] is a programming environment that makes possible the development of parallel and distributed applications. It offers a coordination language to express parallel programs in two main parts: a module graph which defines how nodes interact using data flow streams, and a set of modules, either sequential or parallel, which actually implement the nodes of the graph; in addition, a module or a whole graph can be wrapped as a component interoperable with Web Services. ASSIST and COMPSs have distinct purposes: while the former gives support to high-performance grid-aware applications, the latter offers a much simpler programming model that is oriented to grid-unaware applications.

With respect to workflow managers, some systems have been proposed to specify the elements of a workflow and the connections between them, either graphically or by means of a high-level workflow description language; in this sense they differ from COMPSs, where the workflow graph is implicitly defined by a concrete execution of an application and built automatically and dynamically at runtime. Taverna [23] is a well-known graphical tool for designing and executing grid workflows. A Taverna workflow is specified by a directed acyclic graph where nodes represent software components. Each edge in the graph denotes a data dependency from an output port of the source node to an input port of the destination node. The nodes of a Taverna workflow can be computations executed in the grid and also Web Services, similarly to COMPSs. Although the official Taverna distribution only includes an SSH adaptor to submit computations to grid resources, some projects have developed plugins to make Taverna work on top of Globus-based middleware as well. P-GRADE [20] is a general-purpose, workflow-oriented, Globus-based grid portal; it offers a high-level, graphical workflow development system and an execution environment for various grids. Triana [27] also permits to describe applications by dragging and dropping their components and connecting them together to build

a workflow graph; like in COMPSs, Triana workflows can access the grid through JavaGAT. ASKALON [24] is an application development and computing environment that makes it possible, through the use of a portal, to create a UML model of a workflow; in a second step, this model is automatically translated to an abstract language that represents the workflow and then given to a set of middleware services for scheduling and reliable execution on the grid. Pegasus [18] is a workflow management system that takes high-level workflow descriptions and automatically maps them to grid resources; Pegasus performs execution site selection, manages the input data and provides directives for data transfer and registration.

## 7 CONCLUSIONS AND FUTURE WORK

This paper has shown how an e-Science application can be easily developed with the COMPSs programming model and then parallelised in heterogeneous grids with the COMPSs runtime. Such application is programmed sequentially, while the user is only responsible for specifying its tasks. No API call or resource management details appears in the application, so that it is portable and agnostic of the underlying infrastructure. All the burden of parallelisation and infrastructure management is left to the COMPSs runtime; this paper has demonstrated how this runtime can deal with grids managed by different middleware, making them interoperable while keeping the application unaware of grid details.

The future work includes supporting the use of logical files in COMPSs executions, possibly by creating a JavaGAT adaptor that manages them; such files are referenced with logical names that can be associated to several physical locations. Furthermore, we plan to extend the locality-aware algorithm to take into account not only the number of input files but also their size when deciding the target resource of a task.

## Acknowledgements

This work has been supported by the following institutions: Universitat Politècnica de Catalunya with a UPC Recerca predoctoral grant; the projects of Computación de Altas Prestaciones V and VI (TIN2007-60625, TIN2012-34557); the Spanish Government with grant SEV-2011-00067 of Severo Ochoa Program. On the other hand, the Ibergrid and the Open Science Grid organisations have granted us access to their infrastructures.

## REFERENCES

- [1] Barcelona Supercomputing Center. <http://www.bsc.es>.
- [2] Cloud Interoperability and Portability Remain Science Fiction. <http://searchcloudcomputing.techtarget.com/feature/Cloud-interoperability-and-portability-remain-science-fiction>.

- [3] Grid Interoperation Now Community Group (GIN-CG). <http://www.ogf.org/gf/groupinfo/view.php?group=gin-cg>.
- [4] Ibergrid 2011 Year Report. [http://www.es-ngi.es/documentos/Ibergrid\\\_report\\\_2011\\\_downloadable.pdf](http://www.es-ngi.es/documentos/Ibergrid\_report\_2011\_downloadable.pdf).
- [5] Iniciativa Nacional Grid. <http://www.gridcomputing.pt>.
- [6] Network File System. <http://www.ietf.org/rfc/rfc3010>.
- [7] Open Grid Forum. <http://www.gridforum.org/>.
- [8] Open Portable Batch System. <http://www.openpbs.org/>.
- [9] Open Science Grid. <http://www.opensciencegrid.org>.
- [10] OSG Document Database. <http://osg-docdb.opensciencegrid.org/>.
- [11] Platform Load Sharing Facility. <http://www.platform.com/workload-management/high-performance-computing>.
- [12] ScalaLife Pilot Applications - DISCRETE. <http://www.scalalife.eu/applications>.
- [13] Spanish National Grid Initiative. <http://www.es-ngi.es/>.
- [14] Virtual Organization Membership Service. <http://edg-wp2.web.cern.ch/edg-wp2/security/voms/>.
- [15] ALDINUCCI, M.—COPPOLA, M.—DANELUTTO, M.—VANNESCHI, M.—ZOCOLLO, C.: ASSIST as a Research Framework for High-Performance Grid Programming Environments. In: Cunha, J. C., Rana, O. F. (Eds.): Grid Computing: Software Environments and Tools, Chapter 10. Springer Verlag, 2006, pp. 230–256.
- [16] ALLEN, G.—DAVIS, K.—GOODALE, T.—HUTANU, A.—KAISER, H.—KIELMANN, T.—MERZKY, A.—VAN NIEUWPOORT, R.—REINEFELD, A.—SCHINTKE, F.—SCHÜTT, T.—SEIDEL, E.—ULLMER, B.: The Grid Application Toolkit: Towards Generic and Easy Application Programming Interfaces for the Grid. Proceedings of the IEEE, Vol. 93, 2005, No. 3, pp. 534–550.
- [17] CARGNELLI, M.—ALLEON, G.—CAPPELLO, F.: OpenWP: Combining Annotation Language and Workflow Environments for Porting Existing Applications on Grids. Proceedings of the 2008 9<sup>th</sup> IEEE/ACM International Conference on Grid Computing (GRID'08), Washington, DC, USA, IEEE Computer Society, 2008, pp. 176–183.
- [18] DEELMAN, E.—SINGH, G.—SU, M.-H.—BLYTHE, J.—GIL, A.—KESSELMAN, C.—MEHTA, G.—VAHI, K.—BERRIMAN, G. B.—GOOD, J.—LAITY, A.—JACOB, J. C.—KATZ, D. S.: Pegasus: A Framework for Mapping Complex Scientific Workflows Onto Distributed Systems. Scientific Programming Journal, Vol. 13, 2005, No. 3, pp. 219–237.
- [19] FOSTER, I.—KESSELMAN, C.: Globus: A Metacomputing Infrastructure Toolkit. International Journal of Supercomputer Applications, Vol. 11, 1997, No. 2, pp. 115–128.
- [20] KACSUK, P.—SIPOS, G.: Multi-Grid, Multi-User Workflows in the P-GRADE Grid Portal. Journal of Grid Computing, Vol. 3, 2005, No. 3-4, pp. 221–238.
- [21] LAURE, E.—GRANDI, C.—FISHER, S.—FROHNER, A.—KUNSZT, P.—KRENEK, A.—MULMO, O.—PACINI, F.—PRELZ, F.—WHITE, J.—BARROSO, M.—BUNCIC, P.—BYROM, R.—CORNWALL, L.—CRAIG, M.—DI MEGLIO, A.—DJAOUI, A.—GIACOMINI, F.—HAHKALA, J.—

- HEMMER, F.—HICKS, S.—EDLUND, A.—MARASCHINI, A.—MIDDLETON, R.—SGARAVATTO, M.—STEENBAKKERS, M.—WALK, J.—WILSON, A.: Programming the Grid with gLite. *Computational Methods in Science and Technology*, 2006.
- [22] MCKENNEY, P. E.: Is Parallel Programming Hard, And, If So, What Can You Do About It? `kernel.org`, Corvallis, OR, USA, 2012. <http://kernel.org/pub/linux/kernel/people/paulmck/perfbook/perfbook.html>.
- [23] MISSIER, P.—SOILAND-REYES, S.—OWEN, S.—TAN, W.—NENADIC, A.—DUNLOP, I.—WILLIAMS, A.—OINN, T.—GOBLE, C.: Taverna, Reloaded. In: Gertz, M., Ludäscher, B. (Eds.): *SSDBM 2010*, Heidelberg, Germany, June 2010, *Scientific and Statistical Database Management, Lecture Notes in Computer Science*, Vol. 6187, 2010, pp. 471–481.
- [24] FAHRINGER, T.—PRODAN, R.—DUAN, R.—HOFER, J.—NADEEM, F.—NERIERI, F.—PODLIPNIG, S.—QIN, J.—SIDDIQUI, M.—TRUONG, H.-L.—VILLAZON, A.—WIECZOREK, M.: *ASKALON: A Development and Grid Computing Environment for Scientific Workflows*. *Workflows for eScience*, Springer London, 2007, pp. 450–471.
- [25] RAFANELL, R.: *Extensió de COMP Superscalar*. *Memória del Projecte Fi de Carrera*, Universitat Autònoma de Barcelona, 2011.
- [26] TANAKA, Y.—NAKADA, H.—SEKIGUCHI, S.—SUZUMURA, T.—MATSUOKA, S.: Ninf-G: A Reference Implementation of RPC-Based Programming Middleware for Grid Computing. *Journal of Grid Computing*, Vol. 1, 2003, No. 1, pp. 41–51.
- [27] TAYLOR, I.—SHIELDS, M.—WANG, I.—HARRISON, A.: Visual Grid Workflow in Triana. *Journal of Grid Computing*, Vol. 3, 2005, No. 3-4, pp. 153–169.
- [28] TEJEDOR, E.—BADIA, R. M.: COMP Superscalar: Bringing GRID Superscalar and GCM Together. *Eighth IEEE International Symposium on Cluster Computing and the Grid (CCGrid'08)*, Lyon, France, May 2008, pp. 185–193.
- [29] TEJEDOR, E.—EJARQUE, J.—LORDAN, F.—RAFANELL, R.—ÁLVAREZ, J.—LEZZI, D.—SIRVENT, R.—BADIA, R. M.: A Cloud-Unaware Programming Model for Easy Development of Composite Services. *3<sup>rd</sup> IEEE International Conference on Cloud Computing Technology and Science (CloudCom'11)*, Athens, Greece, November 2011.
- [30] TEJEDOR, E.—FARRERAS, M.—GROVE, D.—BADIA, R. M.—ALMASI, G.—LABARTA, J.: A High-Productivity Task-Based Programming Model for Clusters. *Concurrency and Computation: Practice and Experience*, Vol. 24, 2012, No. 18, pp. 2421–2448.
- [31] THAIN, D.—TANNENBAUM, T.—LIVNY, M.: Condor and the Grid. In: Berman, F., Fox, G., Hey, T. (Eds.): *Grid Computing: Making the Global Infrastructure a Reality*. John Wiley & Sons Inc., December 2002.
- [32] VAN NIEUWPOORT, R. V.—WRZESIŃSKA, G.—JACOBS, C. J.—BAL, H. E.: Satin: A High-Level and Efficient Grid Programming Model. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, Vol. 32, 2010, No. 3, Art. No. 9.



**Enric TEJEDOR** received his Ph.D. from the Technical University of Catalonia (UPC, Spain) in 2013. He conducted his doctorate research as a member of the Grid Computing and Clusters group of the Barcelona Supercomputing Center, where he participated in several EU research projects. As part of his Ph.D., he also carried out two internships at the IBM T. J. Watson Research Center (NY, USA). In 2015 he joined the CERN EP-SFT group, where he currently works on a cloud analysis service for physicists and on parallelization of physics software.



**Javier ÁLVAREZ** holds his M.Sc. degree in computer architecture from the Universitat Politècnica de Catalunya (2013), and since 2014 he has been a Ph.D. student at The University of Adelaide. Previously, he was involved in several research projects as a Research Support Engineer in the Grid Computing and Clusters group at the Barcelona Supercomputing Center.



**Rosa M. BADIA** holds her Ph.D. degree from the Universitat Politècnica de Catalunya (1994). Before, she graduated on computer science at Facultat d'Informàtica de Barcelona (UPC, 1989). She lectured and did research at the Computer Architecture Department (DAC) at UPC from 1989 to 2008, where she held an Associate Professor position from 1997 to 2008. She is the manager of the Workflows and Distributed Computing group at the Barcelona Supercomputing Center. She is a Scientific Researcher at the Spanish National Research Council (CSIC).