

A Service-Oriented Architecture for Scientific Computing on Cloud Infrastructures [★]

Germán Moltó¹, Amanda Calatrava¹, and Vicente Hernández¹

Instituto de Instrumentación para Imagen Molecular (I3M). Centro mixto CSIC
Universitat Politècnica de València CIEMAT, camino de Vera s/n, 46022 Valencia,
España

{gmolto,vhernand}@dsic.upv.es, amcaar@ei.upv.es

Abstract. This paper describes a service-oriented architecture that eases the process of scientific application deployment and execution in IaaS Clouds, with a focus on High Throughput Computing applications. The system integrates i) a catalogue and repository of Virtual Machine Images, ii) an application deployment and configuration tool, iii) a meta-scheduler for job execution management and monitoring. The developed system significantly reduces the time required to port a scientific application to these computational environments. This is exemplified by a case study with a computationally intensive protein design application on both a private Cloud and a hybrid three-level infrastructure (Grid, private and public Cloud).

Topics Parallel and Distributed Computing.

1 Introduction

With the advent of virtualization techniques, Virtual Machines (VM) represent a key technology to provide the appropriate execution environment for scientific applications. They are able to integrate the precise hardware configuration, operating system version, libraries, runtime environments, databases and the application itself in a Virtual Machine Image (VMI) which can be instantiated into one or several runnable entities commonly known as Virtual Appliances. With this approach, the hardware infrastructure is decoupled from the applications, which are completely encapsulated and self-contained. This has paved the way for Cloud computing [1, 2], which enables to dynamically provision and release computational resources on demand.

The efficient and coordinated execution of scientific applications on Cloud infrastructures requires, at least: (i) the dynamic provision and release of computational resources (ii) the configuration of VMs to offer the appropriate execution environment required by the applications and (iii) the allocation and

[★] The authors wish to thank the financial support received from the Generalitat Valenciana for the project GV/2012/076 and to the Ministerio de Economía y Competitividad for the project CodeCloud (TIN2010-17804).

execution of the jobs in the virtualized computational resources. This requires the coordination of different Cloud-enabling technologies in order to automate the workflow required to execute scientific application jobs on the Cloud. To that end, we envision a system where users express their application requirements via declarative procedures and the burden of its deployment, execution and monitoring on an IaaS Cloud infrastructure is automated. There are previous studies that aim at using Cloud computing for scientific computing [3, 4]. However, as far as the authors are aware, there is currently no generic platform that provides automated deployment of scientific applications on IaaS Clouds which deals with VMI management, configuration of VMs and the meta-scheduling of jobs to the virtual computing resources. This represents the whole life cycle of scientific application execution on the Cloud.

For that, the main contribution of this paper is to present a service-oriented architecture integrated by the following developed components: i) a generic catalogue and repository system that indexes VMIs together with the appropriate metadata describing its contents (operating system, capabilities and applications), ii) a contextualization system that allows to deploy scientific applications together with its dependences, iii) a meta-scheduler to manage and monitor the execution of jobs inside VMs and to access the generated output data of the jobs with support for computational steering. The usage of such a system would significantly reduce the time required to migrate an application to be executed on the Cloud. The integration of the different components of the architecture enables to abstract many of the details that arise when interacting with Cloud platforms. This would reduce the entry barrier to incorporate the Cloud as a new source of computational power for scientific applications. This way, scientists would focus on the definition of the jobs and delegate on the proposed platform the orchestration of the components to execute the jobs on the provisioned virtualized infrastructure on the Cloud.

The remainder of the paper is structured as follows. First, section 2 introduces the architecture and details the features of the principal components. Later, section 3 addresses a case study for the execution of a protein design scientific application. Finally, section 4 summarises the paper and points to future work.

2 Architecture for Scientific Application Execution on the Cloud

Many scientific applications require the execution of batch jobs, where each job basically consists of an executable file that processes some input files (or command line arguments) and produces a set of files (or data to the standard output) without the user intervention. This is the case of many parameter sweep studies and Bag of Tasks (BoT) applications commonly found in High Throughput Computing (HTC) approaches, where the jobs share common requirements. For these applications, the benefits of the Cloud are two-fold. Firstly, computational resources can be provisioned on demand according to the number of jobs to be executed (and the budget of the user in the case of a public Cloud). Secondly, the

provisioned VMs can be configured for the precise hardware and software configuration required by the jobs. This means that VMs can be reused to perform the execution of multiple jobs.

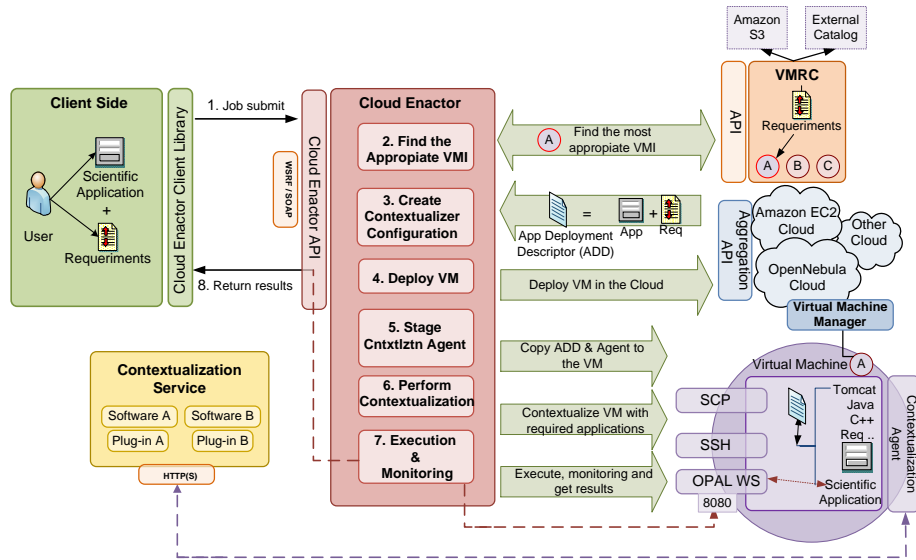


Fig. 1. Interaction diagram for scientific applications execution on IaaS Cloud via the Cloud Enactor

Figure 1 summarises the main interactions between a user and the proposed architecture. The user employs the client-side API to describe each task to be executed (executable file or source code, and required input files) together with the hardware (i.e. CPU architecture, RAM, etc.) and software requirements (OS, applications, system packages, etc.). The jobs might optionally include budget information, since the underlying Cloud infrastructure could require a pay-per-use access to resources. These jobs are submitted (step 1) to the Cloud Enactor (CE) which is the central manager that orchestrates all the components.

The CE checks whether the job could be executed on one of the already deployed (if any) VMs. For the jobs that cannot be executed on the currently deployed VMs, the CE queries the Virtual Machine image & Repository Catalogue (VMRC) [5] with the job's requirements to find the most appropriate VMI to execute the application (step 2). The VMRC, a software that we previously developed, implements matchmaking capabilities to offer a ranked list of suitable VMIs to the Cloud Enactor. The VMRC discards the VMIs that do not satisfy the mandatory requirements (i.e., different OS or CPU architecture) and it ranks the resulting VMIs according to the degree of satisfaction with respect

to the optional requirements (mainly, software applications). The CE computes the deviation from the current state of the most appropriate VMI found and the desired state for the job execution in order to create the Application Deployment Descriptor (ADD) for the contextualization software (step 3). The ADD specifies the deployment process of the application so that the contextualization software can unattendedly perform the installation of the application and its software dependences. This will be executed inside the VM at boot time to deploy the application and its dependences.

Next, the CE must decide the deployment strategy of VMs, which will be in charge of executing the jobs. For that, it has to consider a mixture of performance, economic and trust models to decide the optimum number of VMs to be deployed, together with their Cloud allocation strategy. The performance model should consider the execution time of the jobs (which can be initially estimated by the user but computed after each execution), the deployment time of the VM in the Cloud infrastructure, the time invested in deploying the software requirements of the job (contextualization) and the application itself, as well as the time invested in data transfer, that is, staging out the generated output data of the application inside the VM. The economical model should consider the budget of the user allocated to the execution of each job (or a set of jobs), and the billing policies of the Cloud provider (i.e. hourly rates, economic time zones, etc.). Finally, the trust model plays an important role on scenarios with multiple Cloud providers (Sky Computing), where reputation and the ability of a provider to systematically fulfill the Service Level Agreement (SLA) must be considered. The trust model would be employed to rank a Cloud provider according to its adherence to SLA and the Quality of the Service it offered along the time, among other possible characteristics. For example, a Cloud provider that systematically violates its own SLA should be ranked lower than a provider that has always fulfilled the terms of conditions. The user would express the precise rank function according to the aforementioned categories, as performed in other meta-scheduling softwares such as GridWay.

Therefore, the CE decides to fire up a new VM (or a group of them). This is achieved by delegating on a Virtual Infrastructure Manager (VIM), which deploys the VM on top of a physical infrastructure (step 4). Notice that the CE could use elasticity rules in order to enlarge or shorten the number of VMs dynamically assigned for the allocation of jobs, depending on the budget and the deadline constraints imposed by the user.

When the VM has booted, the CE stages the contextualization agent and the ADD into the VM using SSH (step 5). The VMRC service stores the login name and the private key (or the password) of an account in the VM as part of the metadata stored for a VMI. Then, the contextualization process is started, where software dependences are retrieved from the Contextualization Service and then installed. Next, the scientific application is deployed and a Web services (WS) wrapper is automatically created and deployed into an application server, which is finally started (step 6). This WS wrapper enables to remotely start

and monitor the application running inside the VM. All this automated process results in a VA fully configured for the execution of the scientific application.

Once the VA is up and running, the meta-scheduler can perform the execution of the jobs inside the VAs (step 7). This involves managing and monitoring the execution of the jobs inside the VM during their lifetime. For efficiency purposes each VM would be in charge of the execution of several jobs. In the case of parameter sweep studies and BoT applications commonly found in HTC approaches, the jobs share common requirements and, therefore, they can be executed in the same contextualized VM. In addition, scientific applications might require a periodical access of the generated output data during their executions, mainly for computational steering purposes. Once the application inside the VM has finished executing, then its output data must be retrieved so that another job (with the same requirements) can execute inside the VM.

After all the executions have been carried out, the VAs can be gracefully shutdown which is achieved by the VIM. Notice that it is possible to catalogue the resulting VMI (after the contextualization process) together with the metadata information concerning the new applications installed. Therefore, this would minimize the contextualization time for subsequent executions of that scientific application, since no additional software should have to be installed. This streamlined orchestration of components enables the user to simply focus on the definition of the jobs and thus delegate to the central manager the underlying details of interacting with the Cloud technologies for computational resource provisioning and scientific application execution.

This Service-Oriented Architecture relies on several interoperable services that can be orchestrated by the Cloud Enactor due to the usage of standard protocols and interfaces (WS, WSRF, XML). Concerning the software employed, we have relied on the GMarte meta-scheduler [6], which provides execution management capabilities of scientific tasks on computational Grid infrastructures. By incorporating the functionality to access Cloud infrastructures in this software we can simultaneously schedule jobs on both Grid and Cloud infrastructures. In fact, once the virtual infrastructure of computational resources has been provisioned, other job dispatchers could be fit within the proposed architecture, such as Condor or GridWay. The WS Wrapper for the application is created by the Opal 2 Toolkit [7], which has been integrated in the lightweight contextualization software that we previously developed. Other tools for software configuration, such as Puppet or Chef could also be employed.

3 Case Study

In order to test the suitability of the Cloud infrastructure as a computational source for scientific applications, two case studies were performed. They involve a scientific application that designs proteins with targeted properties via a computationally intensive process based on Monte Carlo Simulated Annealing (MCSA) [8]. The application is developed in the C programming language and it depends

on common build tools available in Linux (configure, make and a C compiler). It also requires the MPICH 2 library, which is a complex software dependence.

For the first case study, we used a fixed number of 8 jobs (an appropriate number for our test infrastructure) and we analysed the total execution time. This time includes from the beginning of the task allocation process until the last job has been executed and its output results have been retrieved. Each job requires the initial configuration of the protein and the matrix that indicates the energetic interactions among the different rotamers of the protein. This amounts to a total of 172 MBytes per job. The job outputs the results of the optimization process to the standard output. This computationally intensive application is typically CPU-bound, but we configured the executions to periodically read the energy matrix from the disk (as part of the optimization process) so that I/O would also be significant in the total runtime.

The test infrastructure is based on four dual-processor Intel Xeon QuadCore with 16 GBytes of RAM Blade servers, with a total of 32 cores, managed by OpenNebula 2.2 and the KVM hypervisor. Two nodes were exclusively used for this particular case study. In order to focus on the execution time, the case study was carried out on pre-started VMs where all the contextualization process had finished and the VMs were ready to receive the execution of the jobs. The allocation of tasks to VMs is achieved by the GMarte meta-scheduler. The current configuration controls that only one job is executed inside a single VM. Therefore, using N VMs allows the concurrent execution of up to N tasks. Other jobs are executed as soon as free VMs are available.

Since the architecture can simultaneously schedule jobs to Grid and both private and public Cloud infrastructures, the second case study executes 30 protein design jobs on a hybrid infrastructure composed by resources from a Grid, the aforementioned private Cloud and the Amazon EC2 public Cloud.

3.1 Results

The solid line in Figure 2.a depicts the global execution time of the case study. As expected, the global execution time decreases when the number of VMs increases, since more computational resources are available to carry out jobs. The plateau in the execution time seen between 4 and 7 VMs is explained by the fact that only one job is executed in each VM and the execution time of each job is expected to be quite similar. Therefore, the executions are actually carried out in groups. As an example, with 5 VMs there is a first group of 5 jobs that are concurrently executed and when they finish, the meta-scheduler allocates the remaining 3 jobs to the free VMs. This would take a similar time as the allocation of the 8 jobs into 7 VMs, which carries out 7 concurrent jobs and a final single job when spare computational resources are available.

The dotted line in Figure 2.a compares the degree of scalability of the Blade servers since it shows the global execution time of the case study when all the VMs are running inside a single node. It can be seen that a similar execution time is achieved except for the case of using 8 VMs, where a minor difference is noticed. Since each node features a dual quad-core processor, it appears that

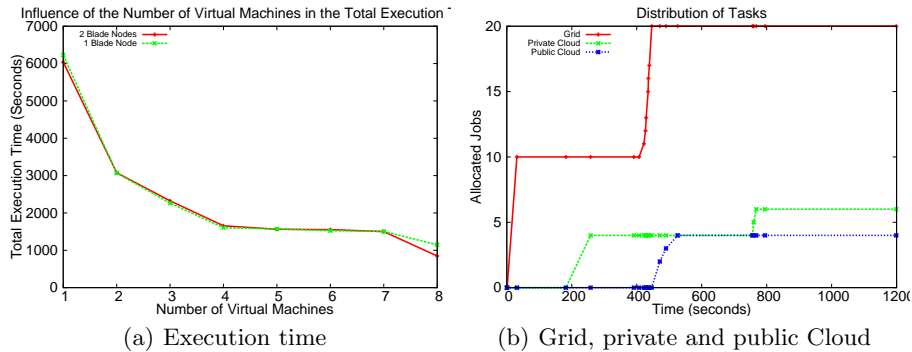


Fig. 2. Global execution time (a) of the case study, considering two different distributions of VMs. Allocated jobs (b) on an infrastructure composed of Grid, private Cloud and public Cloud resources.

scalability issues are only noticeable starting from the 8-th VM in a single node, where the usage of shared resources such as memory and disk start affecting the execution of the applications. These results suggest that VM consolidation in few physical nodes might still deliver good performances for computationally intensive applications, depending on resource consumption.

Concerning the performance improvement gained using the Cloud infrastructure, the results show that up to a speed up of 7.13 is achieved with 8 VMs evenly distributed among the two physical nodes. The global execution time of the case study reduces from a total 6041 seconds in a single VM to just 847 seconds using the aforementioned 8 VMs. Therefore, the usage of virtualised resources from a Cloud as a provider of computational power can deliver a significant improvement for resource-starved scientific applications.

For the second case study, the task allocation is shown on Figure 2.b, further detailed in [9]. The system delegates jobs to a Grid and when that infrastructure is unable to process additional jobs, they are delegated to a provisioned virtual infrastructure from the private Cloud. When neither the Grid or the private Cloud deployment are able to execute more jobs (because all the execution slots are in use), the system can provision computational resources from Amazon EC2 on a pay-per-use basis so that more jobs can be concurrently executed.

To assess that usage pattern we used 10 Grid nodes (from a local resource integrated in the Spanish National Grid Initiative) and 4 provisioned VMs from both the private and the public Cloud. The provisioned VMs were contextualized at boot time in order to deploy the application. We used the Free Usage Tier provided by Amazon EC2 to provision low-performance VMs, thus requiring a noticeably larger time to execute the jobs. Therefore, the developed system allows to simultaneously harvest computational power from three different infrastructures, in order to reduce the execution time of HTC-based applications.

4 Conclusion

This paper has introduced a software architecture that abstracts the details of application deployment and execution on IaaS Clouds. The system features the provision of computational virtualized resources, the configuration of these resources to support the execution of the applications, the cataloguing of virtual machine images and, finally, the job execution management on the virtual infrastructure. The benefits of the proposed architecture have been exemplified by the execution of a protein design case study on both a private Cloud infrastructure and a hybrid infrastructure (Grid, private and public Cloud). The automated deployment and execution of scientific applications fosters the widespread adoption of Cloud technologies by the scientific community. This way, Clouds deliver important benefits for scientific computing in terms of the ability to rapidly provision computational resources and the customizability of the execution environments.

Therefore, the main contribution of this work to the state-of-the-art is the development of generic components and an architecture to integrate them all in order to ease the process of executing scientific applications on the Cloud. Some of the components of the architecture, such as the VMRC system, have been released to the community.

References

1. Vaquero, L.M., Rodero-Merino, L., Caceres, J., Lindner, M.: A break in the clouds. *ACM SIGCOMM Computer Communication Review* **39**(1) (2008) 50
2. Armbrust, M., Fox, A., Griffith, R., Joseph, A.: Above the clouds: A Berkeley view of cloud computing. Technical report, UC Berkeley Reliable Adaptive Distributed Systems Laboratory (2009)
3. Rehr, J., Vila, F., Gardner, J., Svec, L., Prange, M.: Scientific computing in the cloud. *Computing in Science* **99** (2010)
4. Keahey, K., Figueiredo, R., Fortes, J., Freeman, T., Tsugawa, M.: Science Clouds: Early Experiences in Cloud Computing for Scientific Applications. In: *Cloud Computing and its Applications*. (2008)
5. Carrión, J.V., Moltó, G., De Alfonso, C., Caballer, M., Hernández, V.: A Generic Catalog and Repository Service for Virtual Machine Images. In: *2nd International ICST Conference on Cloud Computing (CloudComp 2010)*. (2010)
6. Moltó, G., Hernández, V., Alonso, J.: A service-oriented WSRF-based architecture for metascheduling on computational Grids. *Future Generation Computer Systems* **24**(4) (2008) 317–328
7. Krishnan, S., Clementi, L., Ren, J., Papadopoulos, P., Li, W.: Design and Evaluation of Opal2: A Toolkit for Scientific Software as a Service. In: *2009 IEEE Congress on Services*. (2009)
8. Moltó, G., Suárez, M., Tortosa, P., Alonso, J.M., Hernández, V., Jaramillo, A.: Protein design based on parallel dimensional reduction. *Journal of chemical information and modeling* **49**(5) (2009) 1261–71
9. Calatrava, A. In: *Use of Grid and Cloud Hybrid Infrastructures for Scientific Computing (M.Sc. Thesis in Spanish)*, Universitat Politècnica de València (2012)