



Contents lists available at ScienceDirect

## Future Generation Computer Systems

journal homepage: [www.elsevier.com/locate/fgcs](http://www.elsevier.com/locate/fgcs)

# GAMESH: A grid architecture for scalable monitoring and enhanced dependable job scheduling

Paolo Bellavista<sup>c</sup>, Marcello Cinque<sup>a</sup>, Antonio Corradi<sup>c</sup>, Luca Foschini<sup>c,\*</sup>, Flavio Frattini<sup>a,b</sup>, Javier Povedano-Molina<sup>d</sup>

<sup>a</sup> Dipartimento di Ingegneria Elettrica e delle Tecnologie dell'Informazione - Università di Napoli Federico II - Naples, Italy

<sup>b</sup> RisLab - Laboratorio di Ricerca e Innovazione per la Sicurezza - Naples, Italy

<sup>c</sup> Dipartimento di Informatica - Scienza e Ingegneria - University of Bologna - Bologna, Italy

<sup>d</sup> Real-Time Innovations - Granada, Spain

## HIGHLIGHTS

- A decentralized architecture for monitoring large-scale and multi-domain grids.
- A decentralized job scheduler for large-scale and multi-domain grid environments.
- A novel troubleshooting algorithm for grid data centers.
- Use of standard technologies (e.g., REST, JSON, and DDS) to avoid vendor “lock-in”.

## ARTICLE INFO

### Article history:

Received 5 May 2016

Received in revised form

12 August 2016

Accepted 23 October 2016

Available online xxxx

### Keywords:

Grid  
Monitoring  
Dependability  
Scalability  
Scheduling  
Fault tolerance  
DDS

## ABSTRACT

Grid computing is a largely adopted paradigm to federate geographically distributed data centers. Due to their size and complexity, grid systems are often affected by failures that may hinder the correct and timely execution of jobs, thus causing a non-negligible waste of computing resources. Despite the relevance of the problem, state-of-the-art management solutions for grid systems usually neglect the identification and handling of failures at runtime. Among the primary goals to be considered, we claim the need for novel approaches capable to achieve the objectives of scalable integration with efficient monitoring solutions and of fitting large and geographically distributed systems, where dynamic and configurable tradeoffs between overhead and targeted granularity are necessary. This paper proposes GAMESH, a Grid Architecture for scalable Monitoring and Enhanced dependable job Scheduling. GAMESH is conceived as a completely distributed and highly efficient management infrastructure, concentrating on two crucial aspects for large-scale and multi-domain grid environments: (i) the scalable dissemination of monitoring data and (ii) the troubleshooting of job execution failures. GAMESH has been implemented and tested in a real deployment encompassing geographically distributed data centers across Europe. Experimental results show that GAMESH (i) enables the collection of measurements of both computing resources and conditions of task scheduling at geographically sparse sites, while imposing a limited overhead on the entire infrastructure, and (ii) provides a failure-aware scheduler able to improve the overall system performance, even in the presence of failures, by coordinating local job schedulers at multiple domains.

© 2016 Elsevier B.V. All rights reserved.

## 1. Introduction

The goal of Grid Computing is to federate data centers from several geographically-distributed locations, in order to

achieve large computing and storage capabilities. This computing paradigm is well known in the scientific community since it enables novel experimentations in physics, genetics, astronomy, and more. A notable example is the World-wide Large Hadron Collider (LHC) Computing Grid that, through a joint effort of several countries and data centers around the world [1], allows the scientific community to store and process the huge information volume generated by the CERN LHC project [2]. According to the

\* Corresponding author.

E-mail address: [luca.foschini@unibo.it](mailto:luca.foschini@unibo.it) (L. Foschini).

<http://dx.doi.org/10.1016/j.future.2016.10.023>

0167-739X/© 2016 Elsevier B.V. All rights reserved.

Grid service model, the whole computation load is split into self-contained jobs that are delivered, scheduled, and run on each data center on specific slices of data.

Clearly, improving system performance, by reducing the response time, and increasing the quality of service, in terms of correct job termination, are main goals for the system management. Nevertheless, only a few Grid systems include adequate tools to enable self-adaptive management for reacting to faults and improve overall system performance [3]. Also, most of adopted techniques, such as data and job replication, focus on just one data center [4]. Only a very few efforts consider the possibility to exploit multiple geographically-distributed data centers to move jobs toward more dependable sites. Finally, the monitoring support, when available, typically considers a data center only, rather than complex multi-domain deployments composed by multiple data centers, by addressing mostly performance while widely neglecting fault monitoring and troubleshooting aspects [5].

This paper tackles these issues in Grid system management by proposing a novel solution that exhibits three original characteristics.

- *First*, it proposes an *integrated highly scalable management solution for multi-domain scenarios* that, by continuously monitoring task status and performance metrics across all involved data centers, goes toward full awareness of Grid conditions so to enable failure-aware job scheduling decisions.
- *Second*, it applies optimized techniques for *multi-domain system monitoring*, so to boost communication performance both in single domains and between multiple domains. Within single domains, we propose a monitoring support based on highly scalable industry-level standard solutions. Between different domains, the main goals are interoperability (such as firewall traversal capabilities) and lightweight data transmission.
- *Third*, it leverages on existing hierarchical job scheduling Grid support extended via *novel failure-aware scheduling criteria to deliver incoming jobs to more suitable data centers*.

The proposed solution has been implemented as an open-source support infrastructure, called Grid Architecture for scalable Monitoring and Enhanced dependable job Scheduling (GAMESH), and challenged in a real distributed deployment. GAMESH is available for the Grid community and its monitoring support outperforms related work in the field in terms of delays and scalability [5].

In particular, the reported experimental results show that GAMESH can achieve good and configurable tradeoffs between the imposed monitoring overhead and the flexibility deriving from the usage of high-level monitoring meta-data, especially if compared with related state-of-the-art solutions currently available in the literature. This is confirmed also in challenging deployment scenarios where there is the need to enable the additional costs associated with the GAMESH option for reliable delivery of monitoring data. We have also run extensive experiments to collect performance indicators of GAMESH in inter-domain scenarios, by deploying it on three data-centers across Europe. These inter-domain results show that GAMESH generally requires less bandwidth, with a more predictable behavior, with regards to the existing approaches. Finally, it is demonstrated, through simulation, that the novel failure-aware job scheduling proposal is able to sensibly improve the resource usage. For instance, in our settings, we show that GAMESH can complete 330,000 more jobs over two years with respect to standard job scheduling, without requiring extra computing power.

Let us clearly note that this paper relevantly extends our earlier work appeared in [6]. With respect to it, this extended version provides additional details about the effective design and implementation of selected and primary GAMESH components (see Section 4). In addition, it reports novel and extensive measurements

in both intra-domain and inter-domain deployments. Moreover, it provides the detailed description of the original Stochastic Reward Network models that have been produced to perform the simulation of the GAMESH failure-aware scheduling solution, thus reinforcing also the more theoretical and modeling contribution of this archival article.

The remainder of the paper is organized as follows. Section 2 overviews resource monitoring and job scheduling in multi-domain Grid deployments. Section 3 introduces the system considered as a case study. GAMESH distributed architecture and internal component design are discussed in Section 4. Section 5 reports collected experimental results, with details on the adopted stochastic models. Finally, the paper is concluded in Section 6.

## 2. Related work

The complexity of Grid infrastructures produces extremely challenging monitoring and dependability open issues, especially in large distributed multi-domain deployment scenarios, which require advanced ad-hoc management solutions. In the following, without pretense of being exhaustive, we overview some relevant works in the main areas of monitoring solutions for large distributed computing infrastructures and of scheduling for Grid systems.

**Monitoring system** requirements for Grid environments are defined by the Global Grid Forum in [7]. Such characteristics are related to low latency, high data rate, minimal measurement overhead, security, and scalability. Nevertheless, none of existing monitoring tools satisfies all of them. In the following, some of them are discussed.

A relevant contribution in this area is MonALISA [8]. It is an agent-based monitoring system that supports automatic discovery and achieves scalability by means of proxy services that aggregate monitoring information collected by monitoring sensors periodically. The main draw-back of this architecture is that it relies on centralized monitoring services for storing monitoring data, thus being vulnerable to bottlenecks and single point of failure problems. In [9,10], the authors present Lattice, a framework for monitoring resources on virtualized environments. The most relevant characteristic of Lattice is the communication model it uses: while most monitoring systems are based on request/reply paradigm, Lattice uses the publish/subscribe paradigm to reduce the traffic on the network. Another interesting work in the same direction is [11], where GRIM (Grid Resource Information Monitoring), is introduced. In this work, the authors propose the use of a push approach to reduce the cost of update messages to monitor resources in Grid environments. Another novel contribution in the field of Grid monitoring is presented in [12]. This work does not address the communication problem itself, but it introduces another interesting concept: the inclusion of Complex Event Processing (CEP) in the architecture. This way, it is possible to access to complex queries about the status of monitored resources in real-time. Ganglia is one of the most popular monitoring systems now available and adopted in the Grid [5]. It proposes a scalable hierarchical architecture of services that provide monitoring inside a single data center (single-domain) and across different data centers (multi-domain). Single-domain monitoring is performed by using a proprietary listen/announce multicast protocol; multi-domain monitoring is based on the publishes of monitoring information in an XML file accessible through a socket. Although it has been proved to be a good solution for monitoring Grid, Ganglia does not support dependability mechanisms for high priority data and node control mechanisms. Finally, the widespread Globus Toolkit for Grid management includes a module called the Monitoring and Discovery Service (MDS) [13]. MDS is not a monitoring system itself, but a Web Service-based integration technology to provide standardized

Download English Version:

<https://daneshyari.com/en/article/4950444>

Download Persian Version:

<https://daneshyari.com/article/4950444>

[Daneshyari.com](https://daneshyari.com)