# Semana da Escola de Engenharia
# October 24 - 27, 2011

## IMPROVING THE SUPPORT FOR CRITICAL APPLICATIONS IN GRID BASED COMPUTING INFRASTRUCTURES

Vítor S. P. Oliveira

Centro de Ciências e Tecnologias de Computação

vspo@di.uminho.pt

### ABSTRACT

The number and damages of natural catastrophes have increased significantly over the last decades. Storms and flooding events occur more frequently and since the seventies Europe has seen insured property losses related to flooding alone increase 7% each year (Oliveira 2008). In order to minimize casualties and other losses in crisis situations Civil Protection (CP) agencies call for adequate computing tools which require large amounts of computing power and data storage. The use of the Grid based public computing infrastructure in emergency management can drastically expand the pool of available resources and avoid the high cost of ownership and the single points of failure associated with dedicated infrastructures. FP6 project Cyclops followed on the GMES and INSPIRE EU initiatives to deliver a framework to deploy CP applications in the EGEE Grid infrastructure, which defined the coordinated sharing of applications, computing, storage and communication resources to the agents involved in European CP (Mazzetti 2008).

But the EGEE Grid was designed as a best-effort computing infrastructure and not to support critical applications. By design sites are not expected to be reliable, as the minimum accepted site reliability of only 75% shows (CERN 2010). Notwithstanding the benefits of introducing Quality of Service (QoS) provisions and contracted Service Level Agreements (SLA) in Grid infrastructures (Djemame 2008) current middleware is not adequate for critical CP applications that require workload survivability assurances, that have real-time requirements, that demand emergency job submission, that are latency sensitive or interactive, that need strict data security policies or that integrate with sensor networks and expert systems (Oliveira 2008).

Our proposal for supporting critical CP applications in unreliable computing infrastructures – such as the EGEE Grid, HPC clusters or the Cloud – tries to meet a user specified dependability level by coordinately executing multiple application replicas in several locations and transferring control between them as necessary to provide a high probability of application survivability and predictable execution times. With adequate techniques that assure that application replicas fail independently (Von Neumann 1954), a site reliability of 75% can be converted to 99% (two nines) availability using only 4 application replicas and to five nines using 9 application replicas. To increase the probability of meeting real-time deadlines the replicas evolution can be monitored and the one with the highest performance at each moment can be selected as primary to bypass delays that affect only some of the replicas.

The main research questions are: i) how to distribute the state of a running including the complete execution environment over a wide-area distributed system in a correct, efficient and fault tolerant manner so that it resists failing components without appearing interrupted; ii) how to release applications from the restraints that locality imposes so they can be efficiently moved to other locations; iii) how to observe the applications and analyze their behavior at runtime in order to assess its progress compared to previous executions and to move them to locations better capable of fulfilling the user specified time deadlines.

The presented research builds on the established work on replication and introduces the problems associated with light-weight replication and migration of running applications – provided by the user as opaque virtual machines (VM) – over high latency networks and the online profiling and phase analysis of those applications in order to estimate application progress.

Several methods have been exploited for replicating virtual machines (Bressoud 1995, Yoshiaki 2008, Cully 2008), in which two instances of an application are executed simultaneously in order to support the failure of one of them. The replication of the execution environment preserves not only the internal machine state but also all connections to the rest of the system

including open network connections. Applications are unaware of replication and there is a single image of the application. Also no application data is lost when the primary replica is replaced by the secondary and no restrictions are imposed on how the applications use files, interact with the user sessions or control sensors. The techniques differ between the fully synchronous approaches in which the nodes execute exactly the same code but only one of the them is connected to the exterior, and the asynchronous approaches in which the migration mechanism of the virtual machine monitor is used to synchronize the memory of a virtual machine from a primary node to a secondary node in a kind of permanent pre-migration state. But current approaches to VM replication have little support for n-way replication on wide-area networks and for lightweight migration, as the data storage must either be shared between replicas or completely replicated to each site, making them fit for LAN deployments only. And we know none that monitors the live execution of VMs to use estimations on their future behavior for placement optimization.

Our hypothesis is that the dependability required for critical applications/VMs can be achieved in an unreliable infrastructure if: i) the application data and execution environment are consistently replicated in as many resources as needed to reach the desired level of reliability; ii) there is a lightweight application migration mechanism that can span multiple networks which transparently transfers the control of the application to other replicas upon request or upon failure of the controlling node; and iii) there is a mechanism to analyze application behavior and detect unexpected delays in application execution progress that can endanger meeting the deadlines, in which case control will be transferred to other replica. The research is structured around three main topics. The first is centered on data state preservation with: 1) a model for reliable resource-oriented computing; 2) a reliable resource storage system designed for application mobility. The second topic is application replication over WANs, which includes: 3) application replication and migration to remote networks; 4) network services for WAN replicated applications. The third topic is behavior analysis and estimation with: 5) application monitoring and progress evaluation; 6) exploring the predictability of regular applications.

A dependable autonomic platform was devised to securely deploy user applications (Oliveira 2011). This platform is composed of multi-threaded autonomous agents that cooperate to build a reliable, secure and self-healing overlay system over which all other activities in the system occur, including maintaining a reliable distributed storage system and executing, migrating and monitoring user applications.

## REFERENCES

Bressoud, T. C., Schneider, F. B., 1995, "Hypervisor-based fault tolerance," in Proc. of the 15th ACM SOSP, New York, USA.

CERN, 2012, EGEE Availability and Reliability Report. https://edms.cern.ch/file/963325/1/EGEE_Apr2010.

Cully, B., Lefebvre, G., Meyer, D., Feeley, M., Hutchinson, M., Warfield, A., 2008, "Remus: high availability via asynchronous virtual machine replication", 5th USENIX Symp. Net. Systems Design and Impl., p.161-174, San Francisco.

Djemame, K., Gourlay, I., Padgett, J., Voss, K., Kao, O., 2008, "Risk Management in Grids", in: Market-Oriented Grid Computing, R. Buyya and Kris.Bubendorfer (Editors), Wiley.

Mazzetti, P., Nativi, S., Angelini, V., Verlato, M., Pina, A., Fiorucci, P., 2008, "A Grid Platform for the European Civil Protection e-Infrastructure: the Forest Fires use scenario" 2nd IBERGRID Conference, Porto.

Oliveira, V., 2008, "Civil Protection Applications in a Grid Supported Environment", , MAP-i Doctoral Programme in Computer Science, UM.

Oliveira V., Pina, A., Rocha, A., 2012, "Running User-provided Virtual Machines in Batch-oriented Computing Clusters", PDP'12 (submitted).

Von Neumann, J., "Probabilistic logics and the synthesis of reliable organisms from unreliable components," Automata Studies, no. 34, pp. 43-99, Princeton Univ. Press, 1954.

Yoshiaki, T., Koji, S., Seiji, K., and Satoshi, M., 2008, "Kemari: Virtual machine synchronization for fault tolerance".

## AUTHOR BIOGRAPHY

**VÍTOR OLIVEIRA** is a PhD student working on the thesis entitled "Civil Protection Applications in a Grid Supported Environment". He is supervised by Prof. António Pina and is financially supported by the FCT / UT Austin | Portugal program.