

# Virtualised e-Learning on the IRMOS Real-time Cloud

Tommaso Cucinotta · Fabio Checconi · George Kousiouris · Kleopatra Konstanteli · Spyridon Gogouvitis · Dimosthenis Kyriazis · Theodora Varvarigou · Alessandro Mazzetti · Zlatko Zlatev · Juri Papay · Michael Boniface · Sören Berger · Dominik Lamp · Thomas Voith · Manuel Stein

the date of receipt and acceptance should be inserted later

**Abstract** Providing proper timeliness guarantees to distributed soft real-time applications in a virtualised infrastructure involves the careful use of various techniques at different levels, ranging from real-time scheduling mechanisms at the virtual-machine hypervisor level and QoS-aware protocols at the network level, to proper design methodologies and tools for stochastic modelling and runtime provisioning of the applications. This paper describes the way these techniques were combined to provide strong quality of service guarantees to interactive soft real-time applications in the Cloud Computing infrastructure that has been developed in the con-

text of the IRMOS European Project. The efficiency of the developed infrastructure is demonstrated by two real interactive e-Learning applications, an e-Learning mobile content delivery application and a virtual world e-Learning application, both of which have been integrated into the IRMOS platform.

**Keywords** Real-time scheduling · virtualised infrastructures · stochastic modelling · e-Learning.

## 1 Introduction

A current trend in the software engineering industry is to rely increasingly on the distributed computing paradigm, which is becoming mainstream and ubiquitous, especially due to the high market penetration of low-cost high-speed Internet connectivity. To this direction, more and more applications are developed in a distributed fashion to be hosted on dedicated infrastructures that can be easily accessed by remote users, whether they are using local workstations, laptops, palm-top devices or mobile phones.

This has given rise to the Cloud Computing model, according to which distributed applications developed by Software-as-a-Service (SaaS) providers, are made available by means of tools offered by Platform-as-a-Service (PaaS) providers, to be deployed in a virtualised manner over the resources made available by Infrastructure-as-a-Service (IaaS) providers. Virtualised computing allows for deploying multiple virtual machines (VMs), hosting multiple Operating Systems and services therein, over the same physical hosts, thus achieving a better server consolidation level. Also, the application of virtualisation techniques at network level, has made possible to migrate the VMs from a physical host to another

---

The research leading to these results has received funding from the European Community's Seventh Framework Programme FP7 under grant agreement n.214777 IRMOS – Interactive Realtime Multimedia Applications on Service Oriented Infrastructures.

Tommaso Cucinotta · Fabio Checconi  
Real-Time Systems Laboratory, Scuola Superiore Sant'Anna,  
Pisa, Italy  
E-mail: {t.cucinotta,f.checconi}@sssup.it

George Kousiouris · Kleopatra Konstanteli · Spyridon Gogouvitis · Dimosthenis Kyriazis · Theodora Varvarigou  
National Technical University of Athens, Greece  
E-mail: {gkousiou,kkonst,spyros,g,dimos,dora}@mail.ntua.gr

Alessandro Mazzetti  
eXact learning solutions S.p.A., Sestri Levante, Italy  
E-mail: mazzetti@planasia.eu

Zlatko Zlatev · Juri Papay · Michael Boniface  
University of Southampton, IT Innovation Centre, Southampton,  
UK  
E-mail: {zdz,jp,mjb}@it-innovation.soton.ac.uk

Sören Berger · Dominik Lamp  
University of Stuttgart, Germany  
E-mail: {soeren.berger,dominik.lamp}@rus.uni-stuttgart.de

Thomas Voith · Manuel Stein  
Alcatel Lucent, Stuttgart, Germany  
E-mail: {thomas.voith,manuel.stein}@alcatel-lucent.com

one, for maintenance or load-balancing reasons, in a seamless manner.

However, as the number of VMs deployed over the same physical resources (e.g., links and CPUs) increases, the level of performance experienced by each VM becomes unstable. Indeed it depends heavily on the overall workload imposed by the other VMs competing for the shared resources. Therefore for a virtualised cloud environment to provide proper accommodation to this increasing number of soft real-time distributed applications, proper CPU and network scheduling technologies, coupled with proper performance modelling and runtime provisioning techniques are needed.

In this paper we present the way these techniques have been combined into one virtualised Cloud Computing service-oriented infrastructure that has been developed in the context of the IRMOS European Project<sup>1</sup>, and we show how these concepts have been practically applied to provide strong performance guarantees to real interactive e-Learning applications.

## 2 Related Work

In this section, related works that have appeared in the research literature are briefly summarised.

The work by Lin and Dinda in [19] presents various similarities with the one presented in this paper. First, an EDF-based scheduling algorithm [21] for Linux is used on the host to schedule Virtual Machines (VMs). Furthermore, an analysis is conducted on the application performance, investigating the effects of scheduling decisions and concurrent virtual machines execution. The analysis is very thorough and interesting, however the major limitation of the work resides in the way low-level scheduling is achieved. In fact, the authors make use of a scheduler built into a proper user-space process (VSched), which exploits POSIX real-time priorities in order to achieve an EDF-based scheduling of VMs, and SIGSTOP/SIGCONT signals for realising optionally hard resource reservations. Such an approach presents high scheduling overheads due to the forcibly increased number of context switches, whilst our scheduler [5] is directly built into the kernel and does not introduce any additional context switch; also, VSched cannot properly react to those situations in which a VM blocks or unblocks, e.g., as due to I/O operations, something that is needed in order to guarantee a proper level of temporal isolation, like done instead in our scheduler by exploiting the CBS algorithm [1].

<sup>1</sup> More information is available at: <http://www.irmosproject.eu>.

Another very interesting work by the same authors is [20], where the users of a virtual machine are given the opportunity to adapt the allocated CPU through a simple interface, based on their experience with the application. The cost of the increase is shown, so that the user may decide on the fly. While it is a very promising approach and would eliminate a vast number of issues with regard to application QoS levels, its main drawback is in cases of workflows. Inside a workflow, a degradation in performance may be due to a bottleneck on various nodes executing a part of it. The user will most likely be unaware of the location of the bottleneck, especially in cases of non experts. Instead, the work by Nathuji et al. [24] focuses on automatic online adaptation of the CPU allocation in order to keep a stable performance of VMs. However, the framework does not treat a VM as a “black-box”, but it needs application-specific metrics in order to run the necessary QoS control loop, going beyond the common IaaS business model.

Gupta et al. investigated on the performance isolation of virtual machines [11], focusing on the exploitation of various scheduling policies available in the Xen hypervisor [6]. Furthermore, Dunlap proposed [9] various enhancements to the Xen credit scheduler in order to face with various issues related to the temporal isolation and fairness among the CPU share dedicated to each VM. However, in this work we focus on the KVM<sup>2</sup> hypervisor, along the lines of other works of ours in which we showed how to provide isolation of compute-intensive [7] and network-intensive [8] VMs. Instead, here we also address the modelling issues related to the deployment of an e-Learning application with proper QoS guarantees.

Shirazi et al. [26] proposed DynBench, a benchmark for infrastructures supporting distributed real time applications. This creates dynamic conditions for the testing of the infrastructures. While promising, this framework is mainly oriented towards investigating the limits of the infrastructure and not towards understanding the application behaviour in relationship to different scheduler configurations.

In [10], Germain et al. present DIANE for Grid-based user level scheduling. However, the focus is on controlling the execution end time of long processing applications, and not on real time interactive ones as done in this paper. The problem of optimum allocation of workflows of virtualised services on a set of physical resources under a stochastic approach has been investigated in [14], in the context of soft real-time interactive applications.

<sup>2</sup> More information at: <http://www.linux-kvm.org>.

In terms of application performance modelling in distributed infrastructures a number of interesting works exist. A code analysing process that allows for the simulation of system performance is described in [12]. It models the application by a parameter-driven Conditional Data Flow Graph (CDFG) and the hardware (HW) architecture by a configurable HW graph. The execution cost of each task block in the application CDFG is modelled by user-configurable parameters, which allows for highly flexible performance estimation. The simulator takes the application CDFG and HW graph as the input and performs a low-level simulation to catch the detailed HW activities. While very promising, it needs the source code in order to provide the CDFG. In our work, we deal with VMs as black boxes, what allows for the deployment of applications where the source code is not available for confidentiality purposes.

Another interesting work is presented by Lee et al. in [18]. The application, whose performance must be measured, is run under a strict reservation of resources in order to determine if the given set of reservation parameters satisfies the time constraints for execution. If this is not the case, then these parameters are altered accordingly. If there is a positive surplus, the resources are decreased and if it is negative they are increased until a satisfying security margin is reached. While assuring high utilisation rates, the main disadvantage of this methodology is that this must be performed for every individual execution with the specific SLA parameters before the actual deployment.

Bekner et al. introduce the Vienna Grid Environment (VGE) [3], a framework for incorporating QoS in Grid applications. It uses a performance model to estimate the response time and a pricing model for determining the price of a job execution. In order to determine whether the client's QoS constraints can be fulfilled, for each QoS parameter a corresponding model has to be in place. However, VGE does not prescribe the actual nature of performance models. It specifies only an abstract interface for performance models, taking as granted that these models will be provided from analytical modelling or historical data. But analytical modelling in general requires a thorough knowledge of the application, in order to deduct the equations that depict its performance. In this work we also use analytical models, however they are coupled with a black box approach for the parts of the application that are not visible to the external world besides the developer.

Other works exist that address QoS assurance in Grid environments focusing on performance prediction [16] and control via service selection [17]. While numerous promising solutions exist to the problem of perfor-

mance analysis of VMs in presence of real-time scheduling, these either are not focused on critical parameters that are necessary for running real time applications on SOIs, or they lack for a proper low-level real-time scheduling infrastructure which is needed for supporting temporal isolation among concurrently running VMs.

### 3 Performance Isolation – The IRMOS/ISONI Way

One of the core components which is being developed in IRMOS is the Intelligent (virtualised) Service-Oriented Networking Infrastructure (ISONI) [27]. It acts as a Cloud Computing<sup>3</sup> IaaS provider for the IRMOS framework and manages a set of physical computing, networking and storage resources available in form of multiple nodes/sites within a provider domain (see 1).

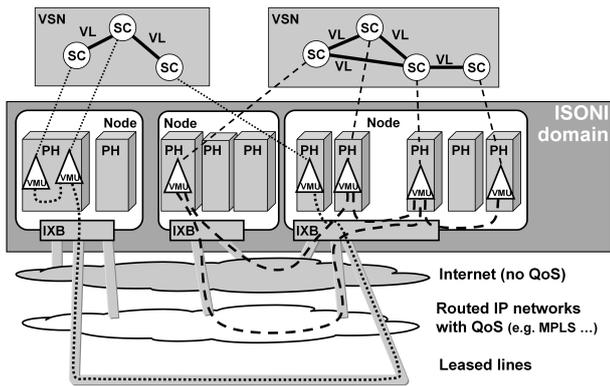
ISONI provides those virtualised resources over which IRMOS applications are deployed. One of the key innovations introduced by ISONI is its capability to ensure guaranteed levels of resource allocation for each individual application instance hosted within the ISONI domain.

This is realised by allowing applications to be deployed in form of a Virtual Service Network (VSN). This is a graph whose vertices represent individual Service Components (SCs) of an application which may be deployed in form of Virtual Machine Units (VMUs), and whose edges represent communications – the virtual links (VLs) – among them.

In order for the system represented by a VSN to comply with real-time constraints as a whole, QoS needs to be supported for all the involved resources, particularly for network links, computing hosts and storage resources. To this purpose, VSN elements are associated with precise resource requirements, e.g., in terms of the required computing performance (e.g. working memory, speed, scheduling) for each node and the required networking performance (e.g. bandwidth, latency, jitter) for each link. These requirements are fulfilled thanks to the allocation and admission control logic pursued by ISONI for instantiating VMs within the managed set of available physical resources, and to the low-level mechanisms shortly described in what follows (a comprehensive ISONI overview is out of the scope of this paper and can be found in [27]).

---

<sup>3</sup> More information at: <http://www.cloudcomputing.org/>.



**Fig. 1** Deployment of Service Components (SCs) within Virtual Machine Units (VMUs) over IRMOS/ISONI.

### 3.1 Isolation of Computing

In order to provide scheduling guarantees to individual VMs scheduled on the same system, processor and core, IRMOS incorporates a hybrid deadline/priority (HDP) real-time scheduler [5] developed within the IRMOS consortium for the Linux kernel. This scheduler provides temporal isolation among multiple possibly complex software components, such as entire VMs (with the KVM hypervisor<sup>4</sup>, a VM is seen as a process). It uses a variation of the Constant Bandwidth Server (CBS) algorithm [1], based on Earliest Deadline First (EDF), for ensuring that each group of processes/threads is scheduled for  $Q$  time units (the budget) every interval of  $P$  time units (the period). The CBS algorithm has been extended for supporting multi-core (and multi-processor) platforms, achieving a partitioned scheduler where the set of tasks belonging to each group may migrate across the associated CBS scheduler instances running on different CPUs, according to the usual multi-processor real-time priority-based scheduling in Linux.

The scheduler exhibits an interface towards user-space applications based on the cgroups [22] framework, which allows for configuration of kernel-level parameters by means of a file system-based interface. This interface has been wrapped within a Python API, in order to make the real-time scheduling services accessible from within the IRMOS platform. The parameters that are exposed by the scheduler are the budget  $Q$  and the period  $P$ , as explained above.

### 3.2 Isolation of Networking

Traffic isolation of independent VSNs within ISONI is achieved by provisioning each VSN deployment with an individual virtual address space and by policing the

network traffic of each deployed virtual link. The automated deployment of policed virtual link overlays avoids unwanted crosstalk between services sharing physical network links and prevents intrusion attempts from unnamed endpoints. The traffic policing avoids that the network traffic traversing the same network elements causes any overload which would lead to an unduly, uncontrolled growth of loss rate, delay and jitter for the network connections of other VSNs. A gap-less policing ensures that the network multiplex stages always get a controlled load of traffic. Therefore, bandwidth policing is an essential building block to ensure QoS for the individual virtual links. It is important to highlight that ISONI allows for the specification of the networking requirements in terms of common and technology-neutral traffic characterisation parameters, such as the needed guaranteed average and peak bandwidth, latency and jitter.

Depending on the specified networking requirements, an adequate transport network is chosen in order to meet the application requirements (see Fig. 1). By default, non-critical traffic without estimable performance requirements can be deployed on Internet transit that does not provide guarantees on the delivery and performance of traffic. On the other hand, soft real-time distributed applications that rely on capacity guarantees can be mapped onto network resources for which ISONI controls the network utilization. The stronger the soft real-time application requirements on delay and jitter become for a virtual link, the shorter is the allowable network distance of the transport resource, which leaves either the site-local network or the use of leased lines as suitable transport resource. Since there exist transport network resources with classification, reservation and other technology-individual mechanisms to enforce QoS of traffic, an ISONI transport network adaptation layer abstracts from transport network technology-individual QoS mechanism like Diffserv [4], Intserv [28][29] and MPLS [25] (see Fig.1).

### 3.3 Performance Modelling

One of the key steps in deploying applications with precise real-time or generally QoS guarantees within IRMOS is the one of building a performance model of the application behavior. This means that, given application-specific configurable parameters (e.g., number of users, resolution of multimedia contents, etc.), and given possible performance levels that one may want to achieve, it should be possible to determine what allocation is needed on the physical resources in order to accomplish that. This is a core information needed

<sup>4</sup> More information at: <http://www.linux-kvm.org>.

by the SaaS provider in order to establish an accurate pricing policy for the customer(s).

Application performance in the cloud depends on many complex factors such as application workload, conditions of the network paths between the user(s) and the server(s) and the computing workload of the physical host(s). Computing workload factors are especially significant in multi-tenant clouds where single hosts are used to service multiple applications. However, the ISONI support for temporal isolation of VMs with guaranteed QoS means the interference due to shared resources becomes negligible. The immediate advantage of this is that the performance of an individual application to be instantiated within the platform may be easily benchmarked and modelled using tools offered by the IRMOS PaaS (see Section 4) as a pure function of its application-specific parameters and the amount of allocated resources, and it turns out to be independent of whatever else is (and/or will be) instantiated in the domain by the provider.

#### 4 The IRMOS PaaS

The IRMOS PaaS, namely the IRMOS Framework Services (FS), acts as a mediator between the SaaS and the ISONI provider. It consists of a family of services and tools that can be used by the SaaS provider for benchmarking, modelling and managing applications on the service-oriented virtualized environment that is offered by the ISONI provider. We can make the following categorization of the IRMOS FS:

- Engineering services: These include services for performing application modelling, benchmarking and performance modelling which are needed by the SaaS provider to make an application ready to be incorporated and hosted in the underlying ISONI cloud.
- Management services: These include services that support the management of the full life-cycle of the service-oriented application. Therefore, the FS include services for discovery, negotiation, reservation, enactment, monitoring and event handling during the execution of an application inside the ISONI virtualized environment.

Using the services above, a description of the VSN can be created and used by the FS for the purpose of negotiating resources with the ISONI provider. This description includes the required resources for the services that make up the application and their interconnections, including instances of the some key FS responsible for real-time critical tasks, such as enactment and monitoring. Upon agreement, the ISONI provider is responsible

for delivering a “just-in-time” deployment of the distributed application on the reserved physical resources according to the specified QoS requirements. From then on, the end user is able to use the application without any knowledge of the underlying Cloud infrastructure.

#### 4.1 Monitoring and Benchmarking

By using the FS tools, the SaaS provider is able to model behavioral aspects of its application, which aid in the performance modelling analysis for deriving critical information about the fluctuation of resource utilization. To this direction, a number of important parameters when benchmarking the application needs to be defined. These include the critical inputs of the services that influence the produced workload, as well as parameters whose values can be chosen by the customer depending on its needs, e.g. the maximum number of clients that are allowed to connect to an e-learning application.

Another factor is the QoS outputs of the services that comprise the application. These are collected and stored during benchmark runs by the FS monitoring system. The latter is a suite of components for monitoring, evaluating possible performance anomalies, and visualizing the performance of the application inside the VSN, along with the levels of the resources that are offered to the application by the ISONI provider. These components acquire this information through an extensible interface that is used to execute application-specific external programs to acquire high-level data about its performance, and also to establish communication with ISONI’s monitoring system for retrieving low-level information about computing and networking resources of the VSN. More information about the design and implementation of the FS monitoring system can be found in [13].

The monitoring information is aggregated and stored in the PaaS infrastructure, and is used both for benchmarking and performance modelling of an application to be deployed in the ISONI. Through this exchange of information, the FS are able to benchmark the application and generate the rules needed for the mapping of high level parameters (used by the SaaS provider for describing the QoS levels), to low-level parameters (used by the ISONI provider for the deployment of the VSN). These rules correlate the configurable workload inputs, the QoS outputs and the hardware assignments of the benchmark runs, and are used in conjunction with analytical approaches to generate an overall end-to-end performance model for the application. This model enables various trade-offs and configurations in order to select the most suitable combination of resources at the

lowest cost, while guaranteeing the QoS levels needed by the SaaS (see Section 3.3).

#### 4.2 Run-time QoS Provisioning

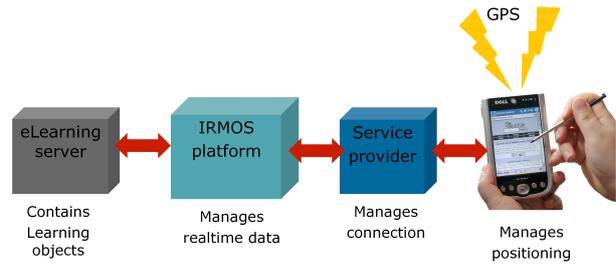
In a Cloud environment, the set of resources available for use can be subjected to frequent changes. For example, new resources (compute servers, file servers, etc) may be added, old ones may be removed or become temporarily unavailable for maintenance purposes, etc. Furthermore, there is always the chance of bad performance modelling that may result in deviations between the agreed and achieved QoS levels at application run-time.

For these reasons, there is a need of continuous observation of the resources and the application's performance for purposes such as fixing problems and tracking down their origin. Because the business entities behind the SaaS, PaaS and IaaS, often have conflicting interests and different mechanisms and levels of fault tolerance, assessment on the source of the problem can be a significant undertaking. To this end, monitoring data are collected during run-time and are being evaluated by the FS monitoring system to identify possible deviations, the source of the problem, and corrective actions that could be undertaken. Typical corrective actions include adjustments of the resource allocation to the application resource utilization, reconfiguration of a running steering service of the application, notifying the end user and the SaaS provider when certain events that affect the availability of the application occur, among others.

Apart from these corrective actions, the FS also enable the SaaS provider to define a range of application-driven events that may automatically trigger certain actions by the FS during the run-time of the application. For example, it may be desirable for the SaaS provider that the FS automatically launch a renegotiation/reconfiguration process when certain usage terms are reached. Such scenarios can be found in a variety of soft real-time applications. In the specific case of a virtual world application such as the one described in Section 5.2, the number of users could be modelled into a trigger for launching automatic renegotiation and reconfiguration of the application by the FS (see Section 8).

### 5 E-Learning Applications

We consider two interactive e-Learning applications with soft real-time requirements:



**Fig. 2** Components of the e-Learning application.

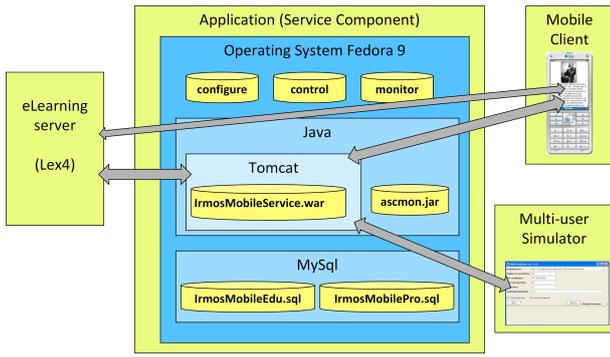
- Mobile Outdoor e-Learning, which is a mobile application, applied to outdoor usage, and
- Virtual World e-Learning, an ubiquitous application for enabling cooperative e-Learning.

#### 5.1 Mobile Outdoor e-Learning

We focus on an e-Learning mobile instant content delivery application, developed to take advantage of a service-oriented architecture paradigm, in which real-time requirements play an important role. In this scenario a user can receive on a mobile phone some e-Learning contents relevant to the current geographical position (e.g., when approaching historical monuments). It consists of a Tomcat<sup>5</sup>-based e-Learning server that exploits a MySQL database for content management (see Fig. 3). The application is able to receive queries with GPS data from multiple clients, search the database and respond with e-Learning contents corresponding to the provided GPS coordinates (see Fig. 2). The application server is provided as a Web Application Archive (war) file, installed on Tomcat, and made available as a Virtual Machine image within the IRMOS infrastructure. Using ISONI, each instance of the application can be assigned precise computing and networking resources to ensure that the high-level requirements defined within an application provider's Service-Level Agreement (SLA) can be met with an agreed level of reliability.

The timing requirements of the application are mainly related to the response times of individual requests submitted by the multitude of users. The response times are gathered by 'ascmon.jar' and communicated to the FS by the 'monitor' script (see Fig. 3). As discussed in Section 3.3, thanks to the deployment within IRMOS, these response times depend merely on high-level application-specific parameters, i.e., on the number of concurrent users querying the same e-Learning instance and the size of the downloaded contents.

<sup>5</sup> More information at: <http://tomcat.apache.org>



**Fig. 3** Software architecture of the Mobile Outdoor e-Learning application.

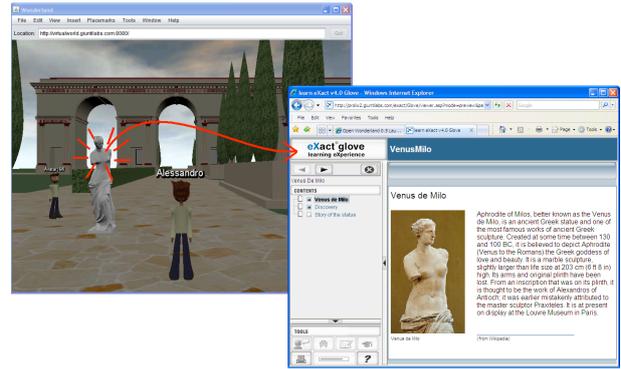
### 5.1.1 Application Client Simulation Description

In order to investigate application performance, we developed a multi-user client simulator. This is capable of simulating the random movements of a certain number of users walking around given GPS coordinates. Then, the simulator mimics the behaviour of the real mobile client associated with the application: whenever the monitored GPS coordinates move sufficiently away from the position of the last queried content, a new request is submitted to the server with the new user position. The number of users and a few parameters governing how each emulated user exactly moves (e.g., the user speed) determine the exact pattern of requests submitted to the server, thus strongly impacting on the imposed server load.

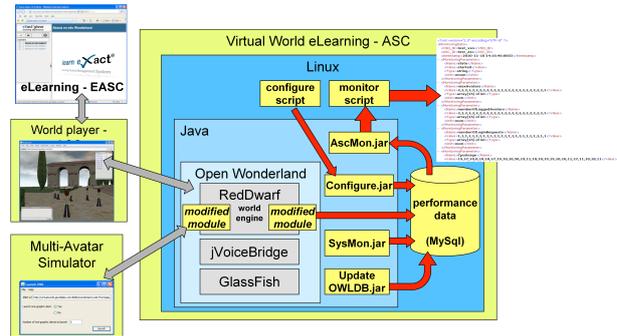
## 5.2 Virtual World e-Learning

In the Virtual World e-Learning application we consider remote users interacting through avatars in a virtual world system. The Virtual World reproduces a museum, where works of art are associated to e-Learning lessons (see Fig. 4). The avatars can move independently from each other and can communicate through a chat or a voice line. When an avatar comes close to a work of art, he or she can invite the other avatars, through chat, to download the corresponding lesson. Each user can play the lesson by themselves, while communicating via chat.

This collaborative virtual world application is composed by two main elements: the application service component (ASC), dealing with user's communication, and the application client component (ACC or World Player), dealing with the graphic rendering. The ASC builds on top of an Open Wonderland server with modified modules for IRMOS adaptation and a MySQL database for managing performance data. A detailed architecture of the Virtual world application is shown in Fig. 5, in which the parts that have been created or



**Fig. 4** The Virtual World application



**Fig. 5** Software architecture of the Virtual World e-Learning application.

adapted to the IRMOS specifications, such as modules for configuration and monitoring, are marked with yellow color. It should also be noted that the Virtual World application provides a multi-avatar simulator for benchmarking purposes. This tool allows the simulation of a large number of users, which are moving continuously to generate the maximum amount of traffic between the clients and the server.

Compared to the Mobile Outdoor e-Learning application, the Virtual World application is more real-time intensive and thus requires a greater number of high level performance parameters such as the avatar speed, the chatting quality, etc. The Virtual World ASC must satisfy real-time constraints (e.g. response time and processing time), in order to support realistically fluid movement of the avatars. The performance of this component strongly depends on the number of connected users, because the quantity of dispatched messages increases polynomially with the number of avatars.

## 6 Performance Estimation

In order to estimate what QoS level can be achieved using different resource configurations, we use performance modelling techniques. Many times, the application internal software structure may be too complex to

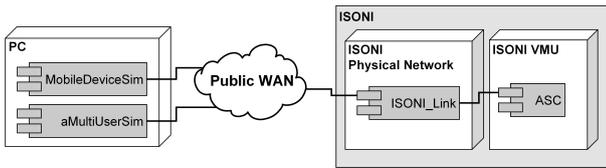


Fig. 6 Modelled elements of the e-Learning application.

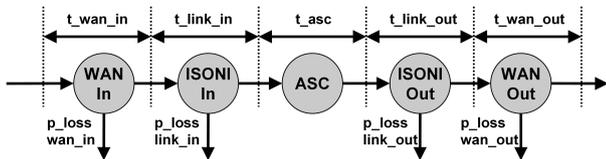


Fig. 7 Stochastic performance model:  $t_{wan\_in}$  and  $t_{wan\_out}$  are modelled as exponential distributions, the other delays as Erlang ones.

be modelled. Or, it may be unknown because developers are reluctant to share detailed information about their application internals, for confidentiality purposes. In other cases, the use of external libraries or components whose internals are unknown makes it impossible to build an exact model. So, from a modelling point of view, it is critical to be able to identify the expected QoS output using a black-box approach.

Therefore, we use a combination of a stochastic model for predicting statistics over the expected run-time networking performance [2], and an Artificial Neural Network (ANN) for identifying the dependency of a component QoS from factors like application-level parameters (e.g., number of clients) or scheduling parameters (e.g., allocated budget and period). These two models, put together, allow for a precise estimate of the overall end-to-end QoS experienced by end-users.

### 6.1 Stochastic Performance Model

We built a Matlab model for simulating, by means of Monte-Carlo type discrete event simulation, a system composed of (see Figure 6): a request generator (modelling the end users), a Public Wide Area Network (WAN), a Private Network internal to an ISONI domain and the VMU hosting the actual Application Service Component (ASC). In order to account for interactivity, we modelled both paths from the user to the ASC and the other way round. The model uses a mix of exponential and Erlang probability distributions (see Figure 7) for modelling the latencies of application requests while traversing the involved networks, and it may also simulate packet loss due to buffer saturation in the various networks (particularly useful for UDP-based communications).

The individual parameters of the model need to be tuned by resorting to proper benchmarking techniques. The behaviour of the latencies inside the ISONI internal network may be accurately estimated thanks to the ISONI networking isolation, and they depend merely on the requested network-level QoS parameters specified in the VSN, and the expected application request pattern. On the other hand, parameters relative to the QoS-unaware WAN must be estimated based on available statistics on the overall network workload foreseen at the time of usage of the application. However, a widespread usage of ISONI would reduce the need for traversing QoS-unaware networks. The behaviour of the ASC was also estimated as an Erlang distribution. However, due to the non-trivial dependency of the performance from application-level parameters, in addition to the resource allocation ones, the Erlang parameters were tuned by resorting to an ANN model (see below).

The described simulator is capable of providing, for each configuration, the full probability distribution of the end-to-end response-times, as well as simpler statistics that may be easily leveraged at design-time, such as the average or a given percentile of the distribution. For example, this allows for finding the configuration parameters granting a given end-to-end response-time with a given probability.

### 6.2 Artificial Neural Networks in the Model

An ANN model is used for modelling the time the server needs in order to retrieve the results from the internal database. The factors that are taken under consideration are mainly the number of connected clients and the scheduling decisions (Q and P parameters). Following the black-box approach, the use of ANNs allows for an easy addition of further inputs (or outputs) as needed (e.g., hardware-specific parameters like the reserved memory or processor speed), once the necessary training data sets are collected.

The investigation of the effect of parameters like the allocated CPU time Q over a period P is critical due to its influence in the QoS output. The choice of P is mainly driven by the time granularity for the allocation needed by the application. For example, for interactive applications, with fast response times and relatively light computations, the granularity must be kept small (in the order of 10–100 ms). For scientific applications performing long and heavy computations, large periods will result in lower overheads (500 ms and beyond).

The outputs of the ANN have been chosen to represent the average and the standard deviation of the expected response times, as due to the configuration

**Table 1** Structure of Mean Response Time and Standard Deviation Prediction ANN.

Layer	Transfer Function		Size (Neurons)
	Mean response time	Standard deviation	
Input	Tansig	Logsig	5
Hidden	Tansig	Logsig	2
Output	Linear (Purelin)	Linear (Purelin)	1
Mean Absolute Error			
Mean response time		2.51%	
Standard deviation		2.75%	

represented at the ANN inputs. These outputs are easily mapped to the Erlang parameters needed for modelling the ASC temporal behaviour in the general model in Figure 7. The ANN model structure is described in V.C., while the data gathered as training set is presented in VI.B.

### 6.3 ANN Structure and Design

In order to implement the ANN, a standard form of network was selected. The type of operation that was desired was function approximation, in order to determine the effect of the input parameters (number of clients, Q, P) on the predicted output (mean value of inner server response time and standard deviation). The collection of the data set was performed with the process described in [15]. Two more inputs were included, CPU speed and VM memory size, but the main focus was on the initial 3 parameters. The resulting network for the mean time prediction was a 3-layer, feed-forward, back propagation network, created through the GNU Octave tool<sup>6</sup>. It was trained with the Octave 'trainlm' function, using the Levenberg-Marquardt algorithm [23]. The structure of the network is shown in Table 1. All the inputs and outputs are normalized in the (-1,1) interval. A standard form of function approximation network was used, with one hidden layer and Tansig transfer functions for the input and hidden layer and linear transfer function for the output layer.

For the standard deviation, a similar process was followed (but with the Logsig transfer function) and the resulting network also appears in Table 1.

## 7 Experimental Results

In this section we report experimental results validating the presented approach to the provisioning of per-

formance guarantees to the Mobile Outdoor e-Learning application by means of proper real-time scheduling and modelling techniques. First, the assumptions of temporal isolation over which the modelling technique relies are validated. Then, some experimental data used for training the ANN models is described, and finally the accuracy of the ANN-based estimations is discussed.

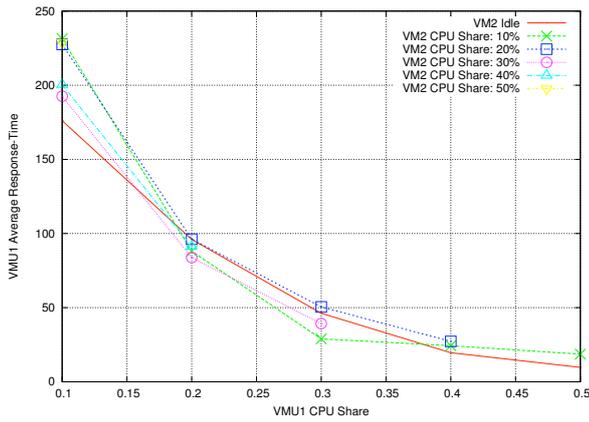
### 7.1 Temporal Isolation by Real-Time Scheduling

We ran an experiment for the purpose of validating our approach to the temporal isolation of VMs concurrently running on the same CPU based on real-time scheduling. To this purpose, we considered two instances of the e-Learning application deployed on the same host and physical core. We launched two instances of the Mobile Outdoor e-Learning multi-user simulator submitting requests coming from 10 emulated users to the two servers, from a second machine in an isolated networking context. We collected the response-times experienced by the two multi-user simulator instances under various conditions in terms of the scheduling parameters configured for the two VMs on the server host.

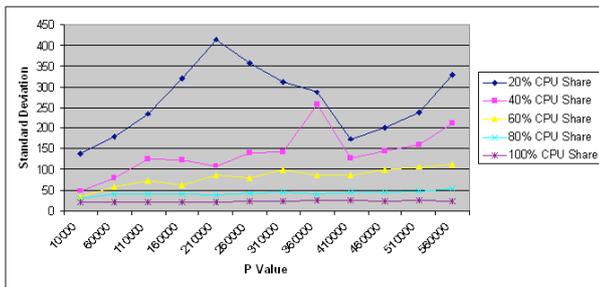
In Figure 8 we report the average response-time of the first VM as a function of the CPU share (on the x-axis) assigned to it, at varying CPU shares assigned to the second competing VM (corresponding to the various curves). Under the ideal conditions of perfect temporal isolation, we would like the second competing VM workload, ranging from completely idle (continuous curve) to having a 50% of load on the system (dashed curve tagged with little triangles), to have no impact at all on the performance of the first VM. This would correspond to having all the curves perfectly superimposed.

As it can be seen from the experimental results, the soft real-time scheduler achieves a nearly good approximation of such a condition, realising a set of curves which are quite close to each other, where the increase of computing resources granted to the second VM corresponds to a slight decrease of the performance of the first VM. This may be mainly attributed to an increased contention on the cache, and constitutes a minimum of interference which cannot be removed. Other factors of interference which are not trivial to keep under control are due to shared resources on the host OS, like the networking stack and interrupts. For example, see [8] for a discussion of the interference due to network-intensive VMs, and the extent to which it can be controlled by real-time scheduling of the CPU.

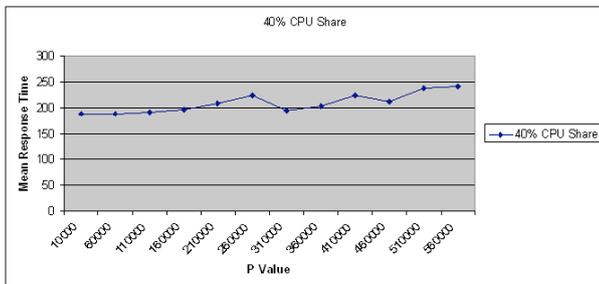
<sup>6</sup> More information is available at: <http://www.gnu.org/software/octave/>.



**Fig. 8** Average response-time of the first VM as a function of the CPU share (on the x-axis) assigned to it, at varying CPU shares assigned to the second competing VM (corresponding to the various curves).



**Fig. 9** Standard Deviation with regard to changing P for 90 users.



**Fig. 10** Mean value with regard to changing P for 70 users and 40% CPU share.

## 7.2 Experimental Performance of the e-Learning Application

In this section, experimental performance data gathered for various configurations of both application-level and resource allocation parameters is presented. The range of values that were altered for the configuration parameters are:

- Number of Users: 30-150
- Q/P (CPU share) : 20-100% with a step of 20
- P: 10000–560000 ( $\mu$ sec) with a step of 50000

For each configuration, about 800 response times were collected, and the corresponding average and standard deviation figures computed. An indicative set of these measurements is discussed below.

The effect of changing granularity on the deviation of the response time values can be observed in Figure 9. This is expected since with high values of P, the service has long active and inactive periods. If the requests fall in the active interval, they will be satisfied quickly but if they fall in the inactive one then they will have to wait until the next active period begins. This effect decreases at increasing allocated CPU shares, since then the CPU is almost dedicated to the application and whenever a request arrives, it is served. The mean response time, as shown in Figure 10, seems not to be affected greatly given that the percentage of CPU assigned is the same.

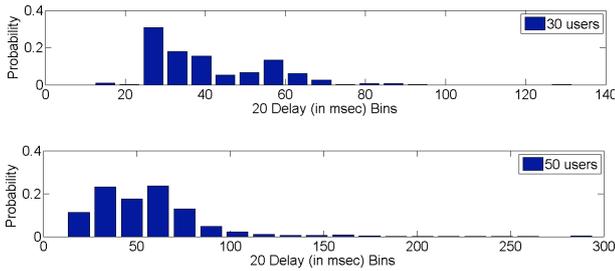
In Figure 11, the comparison of the normalised probability distribution of delay times is shown for two different numbers of users (30 and 50 users). The difference especially in the maximum values of the distributions depicts the effect of the application workload in the response times. Figure 12 shows the cumulative distribution function and the confidence intervals for the same combination of number of users, given that the CPU share allocated to the VMU remains the same.

In Figure 13 all the different configurations are shown for two different numbers of users. In this case, each group of columns (the first high one followed by 4 lower ones) represents one period configuration (P) for different CPU share percentages. The upper line is for low utilization. While the utilization increases the response time decreases. In the horizontal axis, the different P configurations represent increasing period values.

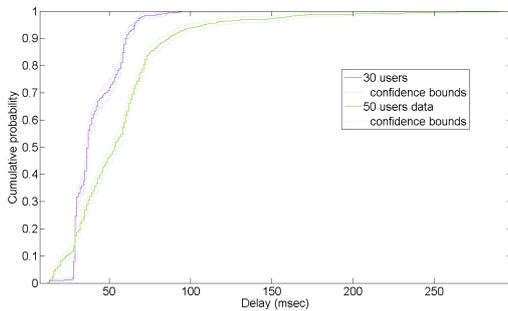
From these measurements it seems interesting that the best granularity (P) should depend also on the percentage of the CPU assigned to the application. In this occasion, for low percentages of utilization it is best to assign values near the middle of the investigated interval (10000-560000 us), as is depicted in Figure 13. For higher percentages of utilisation, lower values of P are more beneficial for the response times of the application. Furthermore, Figure 13 highlights the effect of the increased CPU share allocation to the response time.

## 7.3 Prediction accuracy of the ANN Model

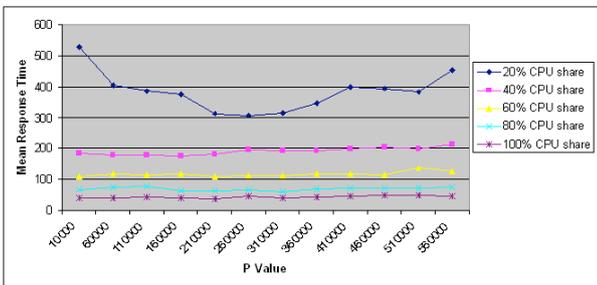
For the estimation of the ANN accuracy, about 30% (87 test executions) of the data set was used only for validation. After the training of the model with the 70% of the test cases, we applied the according inputs of the validation cases and compared the estimated output



**Fig. 11** Comparison of normalised probability distribution of delay times for different number of users.



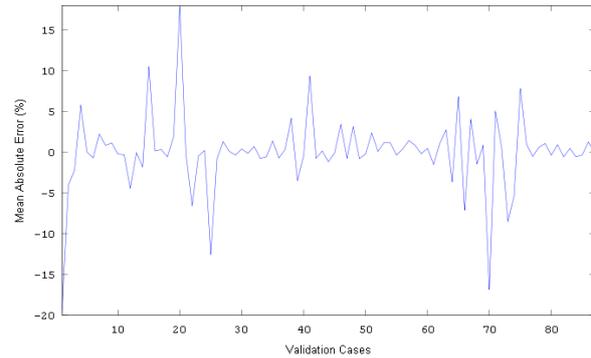
**Fig. 12** Cumulative distribution function and confidence interval for different number of users.



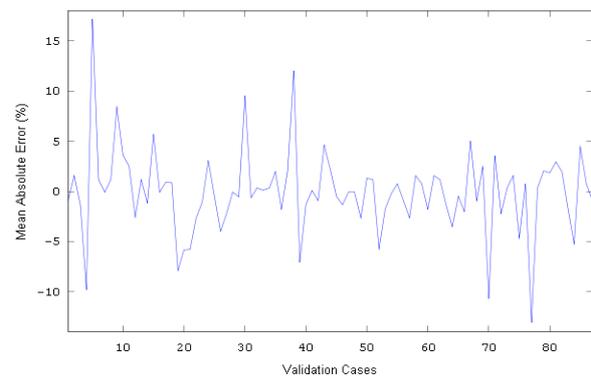
**Fig. 13** Different P's and CPU shares for 110 users.

with the observed one. The overall accuracy was around 2.5% and the error of the network for each individual test case appears in Figure 14. For each validation case, the network error appears in Figure 15.

The accuracy of the ANN models is evident from these measurements, giving sufficient reliability for this part of the overall model. The maximum deviation from the validation cases is very satisfying and so is the mean error for all the experiments. Furthermore, the predictions are not biased, a factor that is critical for the merging of different modelling approaches like in this paper.



**Fig. 14** Accuracy of the ANN for Mean Response Time Prediction.



**Fig. 15** Accuracy of the ANN for Standard Deviation Prediction.

## 8 The Virtual World Use Case

Within the IRMOS context, the Virtual World e-Learning application, as presented in Section 5.2, was used to promote collaborative education through an immersive experience. In more detail, several people (learners, teachers, tutors, organizers, etc.) were allowed to meet in an interactive three-dimensional virtual world environment and interact with 3D-objects that were linked to e-Learning contents (Fig. 4).

This particular use case was used to demonstrate the run-time QoS adaption capabilities behind the IRMOS platform as described in Section 4.2. The negotiation of the Virtual World application is based on a “pay-as-you-reserve” model, i.e. the cost fluctuates according to the amount of resources that are reserved to achieve the requested performance which heavily depends on maximum number of avatars. Therefore in order to maintain the given QoS level, any extra requests for entrance should be rejected when the number of avatars already in the Virtual World is maximum.

However, the IRMOS ISONI is able to support the seamless extension of the resources, by either extending the reserved CPU, memory, network bandwidth share

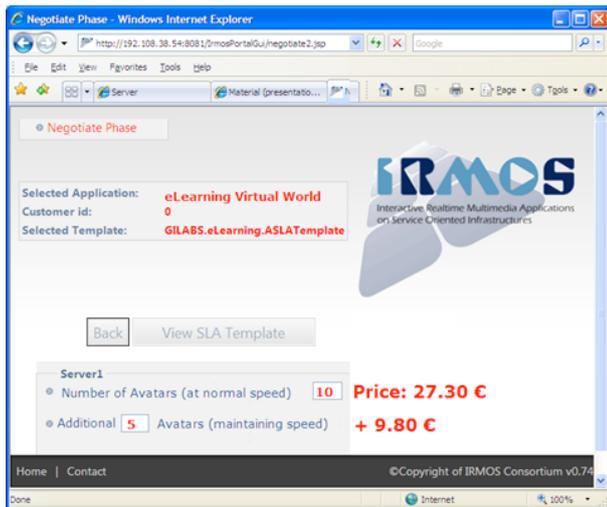


Fig. 16 Negotiation of automatic resource extension.

of the reserved resources or by performing other actions such as live-migrating the application onto other physical resources. Furthermore, the IRMOS FS Monitoring system, as already described in Section 4, is continuously collecting high-level monitoring data during run-time, such as the number of logged avatars and requests for entrance, and can combine them into rules that may trigger certain actions. To this end, the IRMOS customer is given the ability to satisfy a demand for extra resources at run-time, in order to maintain the overall performance without interruption. This has been concretized by a structured negotiation in which the customer can enable the IRMOS FS to automatically re-negotiate an extension of the resources by pre-agreeing to pay an additional cost. Fig. 16 shows the application negotiation in which 5 additional avatars are foreseen.

At the end of the negotiation, the resources for the basic number of avatars (i.e. 10 avatars) are reserved and the application is deployed. Thus up to 10 avatars can enter the Virtual World with the guarantee that the requested QoS level is respected. With 10 already logged avatars, any request after that is being blocked, while the IRMOS FS are re-negotiating with the ISONI provider for resources that are able to sustain the 15 avatars under the same QoS level. In the background, the allocated resources (CPU share, memory, network bandwidth) are being increased according to the output of the performance model for the new maximum number of avatars. At the end of the re-negotiation, a VSN able to support additional users (15 avatars in total) is instantiated and the Virtual World application is reconfigured to accept 15 avatars. The blocked avatar requesting access is now allowed to enter, and the customer is charged with the agreed additional fee.

```

<MonitoringParameter>
  <Name> maxAvatars </Name>
  <Value>10, 10, 10, 10, 10, 10, 10, 15, 15, 15</Value>
  ...
<MonitoringParameter>
  <Name> numberOfLoggedAvatars </Name>
  <Value> 8, 9, 9, 10, 10, 10, 10, 10, 11, 12</Value>
  ...
<MonitoringParameter>
  <Name> numberOfLoginRequests </Name>
  <Value> 3, 2, 2, 1, 1, 2, 2, 2, 1, 0</Value>
  ...
<MonitoringParameter>
  <Name> CpuUsage </Name>
  <Value>40, 40, 40, 40, 40, 60, 60, 60, 60, 60</Value>
  ...

```

Fig. 17 Illustrative example of the QoS adaption process using snapshots (enclosed in a box). Moving left to right: at the first snapshot the renegotiation is triggered, at the second snapshot the reserved CPU share is increased, and at the last snapshot the Virtual World application is reconfigured to accept 15 avatars. The blocked avatars are allowed to log-in afterwards.

The entire process of renegotiation and reconfiguration, including the addition of the blocked avatar in the Virtual World is complete in less than a minute ( $\approx 50$  secs) with the monitoring polling rate set to 10 secs for keeping the monitoring overhead absolutely negligible. However, if appropriate, one can increase the monitoring rate, to achieve a more responsive infrastructure at the cost of higher monitoring overheads imposed on computing and networking resources.

Fig. 17 is an illustrative example of the renegotiation process showing also the XML-based syntax of the gathered monitoring data (see yellow box labeled 'monitor script' in Fig. 5). The monitoring parameters have 10 values, with each of them collected 10 seconds after its previous one.

## 9 Conclusions

In this work, we discussed how two real e-Learning distributed applications have been deployed with predictable and stable QoS levels within the IRMOS platform. We showed how we flanked the temporal isolation mechanisms available within the platform with proper performance analysis, modelling and benchmarking techniques, in order to investigate the performance levels achievable by the application under the various possible configurations. Furthermore, we demonstrated the run-time QoS adaption capabilities behind the IRMOS platform and the way they build proper monitoring and evaluation of application performance on top. In the future, we plan to leverage the black-box approach for performance estimation, so as to apply the described technique to other applications that are already being adapted for deployment within IRMOS, such as distributed editing of professional-quality video and a vir-

tual reality application. Also, we plan to extend the used performance models by accounting for possibly heterogeneous hardware within an ISONI domain. Finally, we plan to extensively compare the predicted performance and the actually realised one, in presence of a variety of other deployed workload types, from compute-intensive to network-intensive ones.

## References

1. L. Abeni and G. Buttazzo. Integrating multimedia applications in hard real-time systems. In *Proceedings of the IEEE Real-Time Systems Symposium*, Madrid, Spain, 1998.
2. Matthew Addis, Zlatko Zlatev, William Mitchell, and Mike Boniface. Modelling interactive real-time applications on service oriented infrastructures. In *2009 NEM Summit - Towards Future Media Internet*, September 2009.
3. S. Benkner and G. Engelbrecht. A generic qos infrastructure for grid web services. In *Proceedings of the Advanced International Conference on Telecommunications and International Conference on Internet and Web Applications and Services (AICT-ICIW)*, 2006.
4. S. Blake, D. Black, M. Carlson, E. Davies, Z. Wang, and W. Weiss. *RFC2475*. An Architecture for Differentiated Service, December 1998.
5. F. Checconi, T. Cucinotta, D. Faggioli, and G. Lipari. Hierarchical multiprocessor CPU reservations for the linux kernel. In *Proceedings of the 5th International Workshop on Operating Systems Platforms for Embedded Real-Time Applications (OSPERT)*, Dublin, Ireland, June 2009.
6. Ludmila Cherkasova, Diwaker Gupta, and Amin Vahdat. Comparison of the three cpu schedulers in xen. *SIGMETRICS Perform. Eval. Rev.*, 35:42–51, September 2007.
7. T. Cucinotta, G. Anastasi, and L. Abeni. Respecting temporal constraints in virtualised services. In *Proceedings of the 2nd IEEE International Workshop on Real-Time Service-Oriented Architecture and Applications (RTSOAA)*, Seattle, Washington, July 2009.
8. T. Cucinotta, D. Giani, D. Faggioli, and F. Checconi. Providing performance guarantees to virtual machines using real-time scheduling. In *to appear in Proceedings of the 5th Workshop on Virtualization and High-Performance Cloud Computing (VHPC)*, Ischia (Naples), Italy, August 2010.
9. G. Dunlap. *Scheduler development update*. Xen Summit Asia 2009, Shanghai, 2009.
10. Cecile Germain-Renaud, Charles Loomis, Jakub Moscicki, and Romain Texier. Scheduling for responsive grids. *Journal of Grid Computing*, 6:15–27, 2008. 10.1007/s10723-007-9086-4.
11. Diwaker Gupta, Ludmila Cherkasova, Rob Gardner, and Amin Vahdat. Enforcing performance isolation across virtual machines in xen. In *Proceedings of the ACM/IFIP/USENIX 2006 International Conference on Middleware*, Middleware '06, pages 342–362, New York, NY, USA, 2006. Springer-Verlag New York, Inc.
12. Zhengting He, Cheng Peng, and Aloysius Mok. A performance estimation tool for video applications. In *Proceedings of the 12th IEEE Real-Time and Embedded Technology and Applications Symposium*, pages 267–276, Washington, DC, USA, 2006. IEEE Computer Society.
13. Gregory Katsaros, George Kousiouris, Spyridon V. Gogouvtis, Dimosthenis Kyriazis, and Theodora A. Varvarigou. A service oriented monitoring framework for soft real-time applications. In *SOCA '10*, pages 1–4, 2010.
14. Kleopatra Konstanteli, Tommaso Cucinotta, and Theodora Varvarigou. Optimum allocation of distributed service workflows with probabilistic real-time guarantees. *Serv. Oriented Comput. Appl.*, 4:68:229–68:243, December 2010.
15. G. Kousiouris, F. Checconi, A. Mazzetti, Z. Zlatev, J. Papay, T. Voith, and D. Kyriazis. Distributed interactive real-time multimedia applications: A sampling and analysis framework. In *Proceedings of the 1st International Workshop on Analysis Tools and Methodologies for Embedded and Real-time Systems (WATERS 2010)*, Brussels, Belgium, July 2010.
16. G. Kousiouris, D. Kyriazis, K. Konstanteli, S. Gogouvtis, G. Katsaros, and T. Varvarigou. A service-oriented framework for gnu octave-based performance prediction. In *Proceedings of the 2010 IEEE International Conference on Services Computing (SCC)*, Miami, Florida, August 2010.
17. Dimosthenis Kyriazis, Konstantinos Tserpes, Andreas Menychtas, Ioannis Sarantidis, and Theodora Varvarigou. Service selection and workflow mapping for Grids: an approach exploiting quality-of-service information. *Concurr. Comput. : Pract. Exper.*, 21:739–766, April 2009.
18. J. W. Lee and K. Asanovic. Meterg: Measurement-based end-to-end performance estimation technique in qos-capable multiprocessors. In *Proceedings of the 12th IEEE Real-Time and Embedded Technology and Applications Symposium (RTAS)*, 2006.
19. B. Lin and P. Dinda. Vsched: Mixing batch and interactive virtual machines using periodic real-time scheduling. In *Proceedings of the IEEE/ACM Conference on Supercomputing*, November 2005.
20. B. Lin and P. Dinda. Towards scheduling virtual machines based on direct user input. In *Proceedings of the 2nd International Workshop on Virtualization Technology in Distributed Computing*, Washington,DC, November 2006.
21. C. L. Liu and James W. Layland. Scheduling algorithms for multiprogramming in a hard-real-time environment. *J. ACM*, 20:46–61, January 1973.
22. P. Menage. *CGROUPS*, 2008. Available on-line at: <http://www.mjmwired.net/kernel/Documentation/cgroups.txt>.
23. Jorge More. The levenberg-marquardt algorithm: Implementation and theory. In G. Watson, editor, *Numerical Analysis*, volume 630 of *Lecture Notes in Mathematics*, pages 105–116. Springer Berlin / Heidelberg, 1978. 10.1007/BFb0067700.
24. R. Nathuji, A. Kansal, and A. Ghaffarkhah. Q-clouds: Managing performance interference effects for qos-aware clouds. In *Proceedings of the 5th European conference on Computer systems (EuroSys)*, Paris, France, April 2010.
25. E. Rosen, A. Viswanathan, and R. Callon. *RFC3031, Multi-protocol Label Switching Architecture*. IETF, January 2001.
26. B. Shirazi, L. Welch, B. Ravindran, C. Cavanaugh, B. Yanamula, R. Brucks, and E. Huh. Dynbench: A dynamic benchmark suite for distributed real-time systems. In *Proceedings of IPDPS Workshop on Embedded HPC Systems and Applications*, S. Juan, Puerto Rico, 1999.
27. T. Voith, M. Kessler, K. Oberle, D. Lamp, A. Cuevas, P. Mandic, and A. Reifert. *ISONI Whitepaper*, September 2008.
28. J. Wroclawski. *RFC2210, The Use of RSVP with IETF Integrated Services*. IETF, September 1997.
29. J. Wroclawski. *RFC2211, Specification of the Controlled Load Quality of Service*. IETF, September 1997.