## Journal of Internet Services and Applications a SpringerOpen Journal

## RESEARCH

**Open Access** 

# Quality-of-service in cloud computing: modeling techniques and their applications

Danilo Ardagna<sup>1</sup>, Giuliano Casale<sup>2\*</sup>, Michele Ciavotta<sup>1</sup>, Juan F Pérez<sup>2</sup> and Weikun Wang<sup>2</sup>

## Abstract

Recent years have seen the massive migration of enterprise applications to the cloud. One of the challenges posed by cloud applications is Quality-of-Service (QoS) management, which is the problem of allocating resources to the application to guarantee a service level along dimensions such as performance, availability and reliability. This paper aims at supporting research in this area by providing a survey of the state of the art of QoS modeling approaches suitable for cloud systems. We also review and classify their early application to some decision-making problems arising in cloud QoS management.

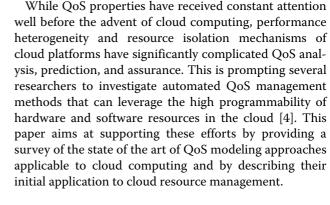
Keywords: Quality of service; Cloud computing; Modeling; QoS management

## 1 Introduction

Cloud computing has grown in popularity in recent years thanks to technical and economical benefits of the ondemand capacity management model [1]. Many cloud operators are now active on the market, providing a rich offering, including Infrastructure-as-a-Service (IaaS), Platform-as-a-Service (PaaS), and Software-as-a-Service (SaaS) solutions [2]. The cloud technology stack has also become mainstream in enterprise data centers, where private and hybrid cloud architectures are increasingly adopted.

Even though the cloud has greatly simplified the capacity provisioning process, it poses several novel challenges in the area of Quality-of-Service (QoS) management. QoS denotes the levels of performance, reliability, and availability offered by an application and by the platform or infrastructure that hosts it<sup>a</sup>. QoS is fundamental for cloud users, who expect providers to deliver the advertised quality characteristics, and for cloud providers, who need to find the right tradeoffs between QoS levels and operational costs. However, finding optimal tradeoff is a difficult decision problem, often exacerbated by the presence of service level agreements (SLAs) specifying QoS targets and economical penalties associated to SLA violations [3].

\*Correspondence: g.casale@imperial.ac.uk



*Scope.* Cloud computing is an operation model that integrates many technological advancements of the last decade such as virtualization, web services, and SLA management for enterprise applications. Characterizing cloud systems thus requires using diverse modeling techniques to cope with such technological heterogeneity. Yet, the QoS modeling literature is extensive, making it difficult to have a comprehensive view of the available techniques and their current applications to cloud computing problems.

*Methodology.* The aim of this survey is to provide an overview of early research works in the cloud QoS modeling space, categorizing contributions according to relevant areas and methods used. Our methodology attempts to maximize coverage of works, as opposed to



© 2014 Ardagna et al.; licensee Springer. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (http://creativecommons.org/licenses/by/4.0), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly credited.

<sup>&</sup>lt;sup>2</sup>Department of Computing, Imperial College London, 180 Queens Gate, London SW7 2AZ, UK

Full list of author information is available at the end of the article

reviewing specific technical challenges or introducing readers to modeling techniques. In particular, we focus on recent modeling works published from 2006 onwards focusing on QoS in cloud systems. We also discuss some techniques originally developed for modeling and dynamic management in enterprise data centers that have been successively applied in the cloud context. Furthermore, the survey considers QoS modeling techniques for interactive cloud services, such as multi-tier applications. Works focusing on batch applications, such as those based on the MapReduce paradigm, are therefore not surveyed.

*Survey Organization.* This survey covers research efforts in *workload modeling, system modeling,* and their *applica-tions* to QoS management in the cloud.

- *Workload modeling* involves the assessment or prediction of the arrival rates of requests and of the demand for resources (e.g., CPU requirements) placed by applications on an infrastructure or platform, and the QoS observed in response to such workloads. We review in Section 2 cloud measurement studies that help characterize those properties for specific cloud systems, followed by a review of workload characterizations and inference techniques that can be applied to QoS analysis.
- *System modeling* aims at evaluating the performance of a cloud system, either at design time or at runtime. Models are used to predict the value of specific QoS metrics such as response time, reliability or availability. We survey in Section 3 formalisms and tools employed for these analyses and their current applications to assess the performance of cloud systems.
- Applications of QoS models often appear in relation to decision-making problems in system management. Techniques to determine optimized decisions range from simple heuristics to nonlinear programming and meta-heuristics. We survey in Section 4 works on decision making for capacity allocation, load balancing, and admission control including research works that provide solutions for the management of a cloud infrastructure (i.e., from the cloud provider perspective) and resource management techniques for the infrastructure user (e.g., an application provider aiming at minimizing operational expenditure, while providing QoS level guarantees to the end users).

Section 5 concludes the paper and summarizes the key findings.

## 2 Cloud workload modeling

The definition of accurate workload models is essential to ensure good predictive capabaility for QoS models. Here, we survey workload characterization studies and related modeling techniques.

## 2.1 Workload characterization

Deployment environment. Several studies have attempted to characterize the QoS showed by cloud deployment environments through benchmarking. Statistical characterizations of empirical data are useful in QoS modeling to quantify risks without the need to conduct an ad-hoc measurement study. They are vital to estimate realistic values for QoS model parameters, e.g., network bandwidth variance, virtual machine (VM) startup times, start failure probabilities. Observations of performance variability have been reported for different types of VM instances [5-7]. Hardware heterogeneity and VM interference are the primary cause for such variability, which is also visible within VMs of the same instance class. Other works characterize the variability in VM startup times [7,8], which is correlated in particular with operating system image size [8]. Some studies on Amazon EC2 have found high-performance contention in CPUbound jobs [9] and network performance overheads [10]. A few characterization studies specific to public and private PaaS hosting solutions also appeared in the literature [11,12], together with comparisons of cloud database and storage services [13-16]. Also, a comparison of different providers on a broad set of metrics is presented in [17].

*Cloud application workloads.* While the above works focus on describing the properties of the cloud deployment environment, users are often faced with the additional problem of describing the characteristics of the workloads processed by a cloud application.

Blackbox forecasting and trend analysis techniques are commonly used to predict web traffic intensity at different timescales. Time series forecasting has been extensively used for web servers for almost two decades. Autoregressive models in particular are quite common in applications and they are already exploited in cloud application modeling, e.g., for auto-scaling [18]. Other common techniques include wavelet-based methods, regression analysis, filtering, Fourier analysis, and kernel-based methods [19].

Recent works in workload modeling that are relevant to cloud computing include [20-22]. Khan et al. [20] uses Hidden Markov Models to capture and predict temporal correlations between workloads of different compute clusters in the cloud. In this paper, the authors propose a method to characterize and predict workloads in cloud environments in order to efficiently provision cloud resources. The authors develop a co-clustering algorithm to find servers that have a similar workload pattern. The pattern is found by studying the performance correlations for applications on different servers. They use hidden Markov models to identify temporal correlations between different clusters and use this information to predict future workload variations. Di et al. [21] defines a Bayesian algorithm for long-term workload prediction and pattern analysis, validating results on data from a Google data center. The authors define nine key features of the workload and use a Bayesian classifier to estimate the posterior probability of each feature. The experiments are based on a large dataset collected from a Google data center with thousands of machines. Gmach et al. [22] applies pattern recognition techniques to data center and cloud workload data. The authors propose a workload demand prediction algorithm based on trend analysis and pattern recognition. This approach aims at finding a way to efficiently use the resource pool to allocate servers to different workloads. The pattern and trend are first analyzed and then synthetic workloads are created to reflect future behaviors of the workload. Zhu and Tung [23] uses a Kalman filter to model the interference caused when deploying applications on virtualized resources. The model accounts for time variations in VM resource usage, and it is used as the basis of a VM consolidation algorithm. The consolidation algorithm is tested and shown to be highly competitive. As the problem of workload modeling is far from trivial, [24] proposes a best practice guide to build empirical models. Important issues are treated, such as the selection of the most relevant data, the modeling technique, and variable-selection procedure. The authors also provide a comparative study that highlights the benefits of different forecasting approaches.

#### 2.2 Workload inference

The ability to quantify resource demands is a pre-requisite to parameterize most QoS models for enterprise applications. Inference is often justified by the overheads of deep monitoring and by the difficulty of tracking execution paths of individual requests [25]. Several works have investigated over the last two decades the problem of estimating, using indirect measurements, the resource demand placed by an application on physical resources, for example CPU requirements. From the perspective of cloud providers and users, inference techniques provide a means to estimate the workload profile of individual VMs running on their infrastructures, taking into account hidden variables due to lack of information.

*Regression Techniques.* A common workload inference approach involves estimating only the *mean* demand placed by a given type of requests on the resource [26-28]. In [26] a standard model calibration technique is introduced. The technique is based on comparing the performance metrics (e.g., response time, throughput and resource utilization) predicted by a performance model against measurements collected in a controlled experimental environment. Given the lack of control over the system workload and configuration during operation, techniques of this type may not be applicable to production systems for online model calibration. These methods exploit queueing theory formulas to relate the mean values of a set of performance metrics (e.g., response times, throughputs, or resource utilizations) to a mean demand to be estimated, e.g., CPU demand. Regression techniques can exploit these formulas to obtain demand estimates from system measurements [29-33].

Zhang et al. [32] presents a queueing network model where each queue represents a tier of a web application, which is parameterized by means of a regression-based approximation of the CPU demand of customer transactions. It is shown that such an approximation is effective for modeling different types of workloads whose transaction mix changes over time.

Liu et al. [33] proposes instead service demand estimation from utilization and end-to-end response times: the problem is formulated as quadratic optimization programs based on queueing formulas; results are in good agreement with experimental data. Variants of these regression methods have been developed to cope with problems such as outliers [34], data multi-collinearity [35], online estimation [36], data aging [37], handling of multiple system configurations [38], and automatic definition of request types [39,40].

Kalbasi et al. [35] proposes the Demand Estimation with Confidence (DEC) approach to overcome the problem of multicollinearity in regression methods. DEC can be iteratively applied to improve the estimation accuracy.

Cremonesi et al. [38] proposes an algorithm to estimate the service demands for different system configurations. A time based linear clustering algorithm is used to identify different linear clusters for each service demands. This approach proves to be robust to noisy data. Extensive validation on generated dataset and real data show the effectiveness of the algorithm.

Cremonesi et al. [39] proposes a method based on clustering to estimate the service time. The authors employ density based clustering to obtain clusters of service times and CPU utilizations, and then use a cluster-wise regression algorithm to estimate the service time. A refinement process is conducted between clustering and regression to get accurate clustering results by removing outliers and merging the clusters that fit the same model. This approach proves to be computationally efficient and robust to outliers.

In [36] an on-line resource demand estimation approach is presented. An evaluation of regression techniques Least Squares (LSQ), Least Absolute Deviations (LAD) and Support Vector Regression (SVR) is presented. Experiments with different workloads show the importance of tuning the parameters, thus the authors proposes an online method to tune the regression parameters. Casale et al. [34] presents an optimization-based inference technique that is formulated as a robust linear regression problem that can be used with both closed and open queueing network performance models. It uses aggregate measurements (i.e., system throughput and utilization of the servers), commonly retrieved from log files, in order to estimate service times.

Pacifici et al. [37] considers the problem of dynamically estimating CPU demands of diverse types of requests using CPU utilization and throughput measurements. The problem is formulated as a multivariate linear regression problem and accounts for multiple effects such as data aging. Also, several works have shown how combining the queueing theoretic formulas used by regression methods with the Kalman filter can enable continuous demand tracking [41,42].

Regression techniques have also been used to correlate the CPU demand placed by a request on multiple servers. For example, linear regression of average utilization measurements against throughput can correctly account for the visit count of requests to each resource [32].

Stepwise linear regression [43] can also be used to identify request flows between application tiers. The knowledge of request flow intensities provides throughputs that can be used in regression techniques.

## 3 System models

Workload modeling techniques presented in Section 2 are agnostic of the logic that governs a cloud system. Explicit modeling of this logic, or part of it, for QoS prediction can help improving the effectiveness of QoS management.

Several classes of models can be used to model QoS in cloud systems. Here we briefly review queueing models, Petri nets, and other specialized formalisms for reliability evaluation. However, several other classes exist such as stochastic process algebras, stochastic activity networks, stochastic reward nets [44], and models evaluated via probabilistic model checking [45]. A comparison of the pros and cons of some popular stochastic formalisms can be found in [46], where the authors highlight the issue that a given method can perform better on some system model but not on others, making it difficult to make absolute recommendations on the best model to use.

## 3.1 Performance models

Among the performance models, we survey queueing systems, queueing networks, and layered queueing networks (LQN). While queueing systems are widely used to model single resources subject to contention, queueing networks are able to capture the interaction among a number of resources and/or applications components. LQNs are used to better model key interaction between application mechanisms, such as finite connection pools, admission control mechanisms, or synchronous request calls. Modeling these feature usually require an in-depth knowledge of the application behavior. On the other hand, while closed-form solutions exist for some classes of queueing systems and queueing networks, the solution of other models, including LQNs, rely on numerical methods.

Queueing Systems. Queueing theory is commonly used in system modeling to describe hardware or software resource contention. Several analytical formulas exist, for example to characterize request mean waiting times, or waiting buffer occupancy probabilities in single queueing systems. In cloud computing, analytical queueing formulas are often integrated in optimization programs, where they are repeatedly evaluated across what-if scenarios. Common analytical formulas involve queues with exponential service and arrival times, with a single server (M/M/1) or with k servers (M/M/k), and queues with generally-distributed service times (M/G/1). Scheduling is often assumed to be first-come first-served (FCFS) or processor sharing (PS). In particular, the M/G/1 PS queue is a common abstraction used to model a CPU and it has been adopted in many cloud studies [47,48], thanks to its simplicity and the suitability to apply the model to multi-class workloads. For instance, an SLA-aware capacity allocation mechanism for cloud applications is derived in [47] using an M/G/1 PS queue as the QoS model. In [48] the authors propose a resource provisioning approach of N-tier cloud web applications by modeling CPU as an M/G/1 PS queue. The M/M/1 open queue with FCFS scheduling has been used [49-51] to pose constraints on the mean response time of a cloud application. Heterogeneity in customer SLAs is handled in [52] with an M/M/k/k priority queue, which is a queue with exponentially distributed inter-arrival times and service times, k servers and no buffer. The authors use this model to investigate rejection probabilities and help dimensioning of cloud data centers. Other works that rely on queueing models to describe cloud resources include [53,54]. The works in [53,54] illustrate the formulation of basic queueing systems in the context of discrete-time control problems for cloud applications, where system properties such as arrival rates can change in time at discrete instants. These works show an example where a non-stationary cloud system is modeled through queueing theory.

*Queueing Networks.* A queueing network can be described as a collection of queues interacting through request arrivals and departures. Each queue represents either a physical resource (e.g., CPU, network bandwidth, etc) or a software buffer (e.g., admission control, or connection pools). Cloud applications are often tiered and queueing networks can capture the interactions between tiers. An example of cloud management solutions exploiting queueing network models is [55], where the cloud service center is modeled as an open queueing network of multiclass single-server queues. PS scheduling is assumed at the resources to model CPU sharing. Each layer of queues represents the collection of applications supporting the execution of requests at each tier of the cloud service center. This model is used to provide performance guarantees when defining resource allocation policies in a cloud platform. Also, [56] uses a queueing network to represent a multi-tier application deployed in a cloud platform, and to derive an SLA-aware resource allocation policy. Each node in the network has exponential processing times and a generalized PS policy to approximate the operating system scheduling.

*Layered Queueing Networks.* Layered queueing networks (LQNs) are an extension of queueing networks to describe layered software architectures. An LQN model of an application can be built automatically from software engineering models expressed using formalisms such as UML or Palladio Component Models (PCM) [57]. Compared to ordinary queueing networks, LQNs provide the ability to describe dependencies arising in a complex workflow of requests and the layering among hardware and software resources that process them. Several evaluation techniques exist for LQNs [58-61].

LQNs have been applied to cloud systems in [62], where the authors explored the impact of the network latency on the system response time for different system deployments. LQNs are here useful to handle the complexity of geo-distributed applications that include both transactional and streaming workloads.

Jung et al. [63] uses an LQN model to predict the performance of the RuBis benchmark application, which is then used as the basis of an optimization algorithm that aims at determining the best replication levels and placement of the application components. While this work is not specific to the cloud, it illustrates the application of LQNs to multi-tier applications that are commonly deployed in such environments.

Bacigalupo et al. [64] investigates a prediction-based cloud resource allocation and management algorithm. LQNs are used to predict the performance of an enterprise application deployed on the cloud with strict SLA requirements based on historical data. The authors also provide a discussion about the pros and cons of LQNs identifying a number of key limitations for their practical use in cloud systems. These include, among others, difficulties in modeling caching, lack of methods to compute percentiles of response times, tradeoff between accuracy and speed. Since then, evaluation techniques for LQNs that allow the computation of response time percentiles have been presented [61]. *Hybrid models*. Queueing models are also used together with machine learning techniques to achieve the benefits of both approaches. Queueing models use the knowledge of the system topology and infrastructure to provide accurate performance predictions. However, a violation of the model assumptions, such as an unforeseen change in the topology, can invalidate the model predictions. Machine learning algorithms, instead, are more robust with respect to dynamic changes of the system. The drawback is that they adopt a black-box approach, ignoring relevant knowledge of the system that could provide valuable insights into its performance.

Desnoyers et al. [43] studies the relations between workload and resource consumption for cloud web applications. Queueing theory is used to model different components of the system and data mining and machine learning approaches ensure dynamic adaptation of the model to work under system fluctuations. The proposed approach is shown to achieve high accuracy for predicting workload and resource usages.

Thereska et al. [65] proposes a robust performance model architecture focusing on analyzing performance anomalies and localizing the potential source of the discrepancies. The performance models are based on queueing-network models abstracted from the system and enhanced by machine learning algorithms to correlate system workload attributes with performance attributes.

A queueing network approach is taken in [66] to provision resources for data-center applications. As the workload mix is observed to fluctuate over time, the queueing model is enhanced with a clustering algorithm that determines the workload mix. The approach is shown to reduce SLA violations due to under-provisioning in applications subject to to non-stationary workloads.

#### 3.2 Dependability models

Petri nets, Reliability Block Diagrams (RBD), and Fault Trees are probably the most widely known and used formalisms for dependability analysis. Petri nets are a flexible and expressive modeling approach, which allows a general interactions between system components, including synchronization of event firing times. They also find large application also in performance analysis.

RBDs and Fault Trees aim at obtaining the overall system reliability from the reliability of the system components. The interactions between the components focus on how the faulty state of one or more components results in the possible failure of another components.

*Petri nets.* It has long been recognized the suitability of Petri nets for performance and dependability of computer systems. Petri nets have been extended to consider stochastic transitions, in stochastic Petri nets (SPNs) and generalized SPNs (GSPNs). They have recently enjoyed a resurgence of interest in service-oriented systems to describe service orchestrations [67].

In the context of cloud computing, we have more application examples of Petri nets nets for dependability assessment, than for performance modeling. Applications to cloud QoS modeling include the use of SPNs to evaluate the dependability of a cloud infrastructure [68], considering both reliability and availability. SPNs provide a convenient way in this setting to represent energy flow and cooling in the infrastructure. Wei et al. [69] proposes the use of GSPNs to evaluate the impact of virtualization mechanisms, such as VM consolidation and live migration, on cloud infrastructure dependability. GSPNs are used to provide fine-grained detail on the inner VM behaviors, such as separation of privileged and non-privileged instructions and successive handling by the VM or the VM monitor. Petri nets are here used in combination with other methods, i.e., Reliability Block Diagrams and Fault Trees, for analyzing mean time to failure (MTTF) and mean time between failures (MTBF).

*Reliability Block Diagrams.* Reliability block diagrams (RBDs) are a popular tool for reliability analysis of complex systems. The system is represented by a set of inter-related blocks, connected by series, parallel, and *k*-out-of-*N* relationships.

In [70], the authors propose a methodology to evaluate data center power infrastructures considering both reliability and cost. RBDs are used to estimate and enforce system reliability. Dantas et al. [71] investigates the benefits of a warm-standby replication mechanism in Eucalyptus cloud computing environments. An RBD is used to evaluate the impact of a redundant cloud architecture on its dependability. A case study shows how the redundant system obtains dependability improvements. Melo et al. [72] uses RBDs to design a rejuvenation mechanism based on live migration, to prevent performance degradation, for a cloud application that has high availability requirements.

*Fault Trees.* Fault Trees are another formalism for reliability analysis. The system is represented as a tree of inter-related components. If a component fails, it assumes the logical value *true*, and the failure propagation can be studied via the tree structure. In cloud computing, Fault Trees have been used to evaluate dependencies of cloud services and their effect on application reliability [73]. Fault Trees and Markov models are used to evaluate the reliability and availability of fault tolerance mechanisms. Jhawar and Piuri [74] uses Fault Trees and Markov models to evaluate the reliability and availability of a cloud system under different deployment contexts. Based on this evaluation, the authors propose an approach to identify the best mechanisms according to user's requirements. Kiran

et al. [75] presents a methodology to identify, mitigate, and monitor risks in cloud resource provisioning. Fault Trees are used to assess the probability of SLA violations.

#### 3.3 Black-box service models

Service models have been used primarily in optimising web service composition [76], but they are now becoming relevant also in the description of SaaS applications, IaaS resource orchestration, and cloud-based business-process execution. The idea behind the methods reviewed in this section is to describe a service in terms of its response time, assuming the lack of any further information concerning its internal characteristics (e.g., contention level from concurrent requests).

Non-parametric blackbox service models include methods based on deterministic or average execution time values [77-81]. Several works instead adopt a description that includes standard deviations [76,82,83] or finite ranges of variability for the execution times [84,85]. Parametric service models instead assume exponential or Markovian distributions [86,87], Pareto distributions to capture heavytailed execution times [88], or general distributions with Laplace transforms [89].

Huang et al. [90] presents a graph-theoretic model for QoSaware service composition in cloud platforms, explicitly handling network virtualization. Here, the authors explore the QoS-aware service provisioning in cloud platforms by explicitly considering virtual network services. A system model is demonstrated to suitably characterize cloud service provisioning behavior and an exact algorithm is proposed to optimize users' experience under QoS requirements. A comparison with state of the art QoS routing algorithms shows that the proposed algorithm is both cost-effective and lightweight.

Klein et al. [91] considers QoS-aware service composition by handling network latencies. The authors present a network model that allows estimating latencies between locations and propose a genetic algorithm to achieve network-aware and QoS-aware service provisioning.

The work in [92] considers cloud service provisioning from the point of view of an end user. An economic model based on discrete Bayesian Networks is presented to characterize end-users long-term behavior. Then the QoS-aware service composition is solved by Influence Diagrams followed by analytical and simulation experiments.

#### 3.4 Simulation models

Several simulation packages exist for cloud system simulation. Many solutions are based on the CLOUDSIM [93] toolkit that allows the user to set up a simulation model that explicitly considers virtualized cloud resources, potentially located in different data centers, as in the case of hybrid deployments. CLOUDANALYST [94] is an extension of CLOUDSIM that allows the modeling of geographically-distributed workloads served by applications deployed on a number of virtualized data centers.

EMUSIM [95] builds on top of CLOUDSIM by adding an emulation step leveraging the Automated Emulation Framework (AEF) [96]. Emulation is used to understand the application behavior, extracting profiling information. This information is then used as input for CLOUDSIM, which provides QoS estimates for a given cloud deployment.

Some other tools have been developed to estimate data center energy consumption. For example, GREEN-CLOUD [97], which is an extension of the packet-level simulator NS2 [98], aims at evaluating the energy consumption of the data center resources where the application has been deployed, considering servers, links, and switches.

Similarly, DCSIM [99] is a data center simulation tool focused on dynamic resource management of IaaS infrastructures. Each host can run several VMs, and has a power model to determine the overall data center power consumption.

GROUDSIM [100] is a simulator for scientific applications deployed on large-scale clouds and grids. The simulator is based on events rather than on processes, making it a scalable solution for highly parallelized applications.

*Research Challenges* A threat to workload inference on IaaS clouds is posed by resource contention by other users, which can systematically result in biased readings of performance metrics. While some bias components can be filtered out (for example using the CPU steal metric available on Amazon EC2 virtual machines), contention on resources such as cache, memory bandwidth, network, or storage, is harder or even impossible to monitor for the final user. Research is needed in this domain to understand the impact of such contention bias on demand estimation.

Major complications arise in workload inference on PaaS clouds, where infrastructure-level metrics such as CPU utilization are normally unavailable to the users. This is a major complication for regression methods which all depend on mean CPU utilization measurements. Methods based on statistical distributions do not require CPU utilization, but they are still in their infancy. More work and validations on PaaS data are required to mature such techniques.

## **4** Applications

A prominent application of QoS models is optimal decision-making for cloud system management. Problem areas covered in this section include capacity allocation,

load balancing, and admission control. Several other relevant decision problems exist in cloud computing, e.g., pricing [101], resource bidding [102], and provider-side energy management [103].

We classify works in three areas using, for comparability, a taxonomy similar to the one appearing in the software engineering survey of Aleti et al. [104], which also covers design-time optimization studies, but does not focus on cloud computing. Our classification dimensions follow from these questions:

*Perspective:* is the study focusing on the perspective of the infrastructure user or or on the perspective of the provider?

*Dimensionality:* is the study optimizing a single or multiple objective functions?

*Solution:* is the presented solution centralized or distributed?

*Strategy:* is the optimization problem tackled by an exact or approximate technique?

*Time-scale:* is the time-scale for the performed adaptations, which can be short (seconds), medium (minutes), or long (hours/days)?

*Discipline:* is the management approach based on control theory, machine learning or operations research (i.e., optimization, game theory, bio-inspired algorithms)?

Table 1 provides a taxonomy of the papers reviewed in the next sections, organized according to the above criteria. Few remarks are needed to clarify the methodology used to classify the papers:

- In the *Perspective* dimension, a public PaaS or SaaS service built on top of a public IaaS offering is classified as a user-side perspective.
- Under *Dimensionality*, we treat studies that weight multiple criteria into a single objective as Multi-Objective methods.

Finally, the following observations on the *Discipline* dimensions must also be made.

- *Control theory* has the advantage of guaranteeing the stability of the system upon workload changes by modeling the transient behavior and adjusting system configurations within a transitory period [143].
- *Machine learning* techniques, instead, use learning mechanisms to capture the behavior of the system without any explicit performance or traffic model and with little built-in system knowledge. Nevertheless, training sessions tend to extend over several hours [144] and retraining is required for evolving workloads.
- *Operations research* approaches are designed with the aim of optimizing the degree of user satisfaction. The goals, in fact, are expressed in terms of user-level QoS

## Table 1 Decision-making in the cloud - a taxonomy

Category	Value	Application		
		Capacity allocation	Admission control	Load balancing
Perspective	Infrastructure user	[47] [105][106] [107] [108] [109] [110] [111] [112] [113] [114] [115] [116] [117]	[118] [116] [119]	[120] [47] [105] [107] [121] [122] [123]
	Infrastructure provider	[55] [124] [49] [125] [126] [52] [56] [127] [128] [129] [130] [131] [54] [132] [18] [133] [134] [135] [136] [137] [138] [51] [23]	[49] [131] [130] [139] [140] [124] [52]	[50] [141] [132] [134]
Dimensionality	Single-Objective	[55] [47] [106] [132] [49] [105] [108] [126] [56] [127] [128] [109] [130] [110] [111] [112] [113] [114] [115] [134] [116] [135] [136] [138] [51] [23] [129]	[49] [116] [130] [119] [118] [139] [140]	[120] [122] [105] [132] [134] [47] [123] [141] [50]
	Multi-Objective	[125] [52] [137] [133] [124] [107] [131] [18] [117]	[52] [124] [131]	[121] [107]
Solution	Centralized	[124] [49] [108] [125] [126] [56] [132] [127] [128] [130] [131] [110] [111] [112] [18] [133] [113] [114] [115] [134] [116] [135] [137] [138] [52] [136] [23] [129] [51] [117]	[124] [49] [142] [52] [131] [116] [130] [119] [118] [139] [140]	[120] [122] [50] [132] [121] [141] [134]
	Decentralized	[47] [105] [106] [109] [107] [55] [54]		[47] [105] [107] [123]
Strategy	Exact	[52] [131] [116] [108] [49] [113] [114] [115] [51]	[52] [49] [142] [131] [116]	[122]
	Approximate	[47] [105] [132] [124] [106] [130] [109] [108] [23] [129] [125] [126] [56] [110] [111] [112] [18] [138] [133] [135] [134] [137] [127] [107] [55] [54] [136] [117]	[124] [130] [119] [118] [139][140]	[50] [120] [105] [132] [47] [107] [123] [121] [141] [134]
Timescale	Short	[47] [105] [52] [128] [54] [113] [115]	[142] [49] [52] [139] [140]	[120] [122] [47] [105] [141]
	Medium	[124] [49] [47] [132] [52] [131] [109] [108] [125] [126] [56] [129] [110] [111] [112] [18] [138] [106] [133] [134] [137] [127] [107] [55] [54] [130] [136] [51] [23] [114] [117]	[124] [131] [52] [130] [119]	[47] [132] [107] [123] [134] [121]
	Long	[116] [108] [111] [135] [55]	[116] [118]	[50]
	Control Theory	[54] [126] [110] [112] [114]		
Discipline	Machine Learning	[110] [134] [23]	[140]	[134]
	Operations Research	[124] [49] [47] [52] [128] [130] [131] [116] [109] [108] [132] [125] [56] [129] [110] [111] [18] [106] [136] [138] [133] [135] [137] [127] [107] [55] [115] [113] [117] [51]	[124] [49] [142] [52] [119] [139] [116] [131] [130] [141]	[120] [122] [132] [47] [107] [123] [121] [50]

metrics. Typically, this approach consists of a performance model embedded within an optimization program, which is solved either globally, locally, or heuristically.

## 4.1 Capacity allocation

## 4.1.1 Infrastructure-provider capacity allocation

The capacity allocation problem arising at the provider side involves deciding the optimal placement of running

applications on a suitable number of VMs, which in turn has to be executed on appropriate physical servers. The rationale is to assign resource shares trying to minimize management costs (formed mainly by costs associated with energy consumption), while guaranteeing the fulfilment of SLAs stipulated with the customers. Bin packing is a common modeling abstraction [51], but its NP-hardness calls for heuristic solutions. In [127] the capacity allocation problem is solved by means of a dynamic algorithm, since static allocation policies and pricing usually lead to inefficient resource sharing, poor utilization, waste of resources and revenue loss when demands and workloads are time varying. The paper presents a Minimum Cost Maximum Flow (MCMF) algorithm and compares it against a modified Bin-Packing formulation; the MCMF algorithm exhibits very good performance and scalability properties. An autoregressive process is used to predict the fluctuating incoming demand.

In [137], autoscaling is modeled as a modified Class Constrained Bin Packing problem. An auto-scaling algorithm is provided that automatically classifies incoming requests. Moreover, it improves the placement of application instances by putting idle machines into standby mode and reducing the number of running instances in condition of light load.

Tang et al. [133] proposes a fast heuristic solution for VM placement over a very large number of servers in a IaaS data center to equally balance the CPU load among physical machines, taking into account also memory requirements of running applications.

In [125] is proposed a framework for VM deployment and reconfiguration optimization, with the aim at increasing profits of IaaS providers. The authors reduce costs considering the balance of multi-dimensional resources utilization and building up an optimization method for resource allocation; as far as reconfiguration is concerned, they propose a strategy for VM adjustment based on time-division multiplex and on VM live migration.

In [55] a VM placement problem for a PaaS is solved at multiple time-scales through a hierarchical optimization framework. Authors in [132] provide a solution for trafficaware VM placement minimizing also network latencies among deployed applications. The work presents a twotier approximate algorithm able to successfully solve very large problem instances. Moreover, a formulation for the considered problem is presented and its hardness is proven.

A capacity allocation problem is also studied in [136], in which a game-theoretic method is used to find approximated solutions for a resource allocation problem in the presence of tasks with multiple dependent subtasks. The initial solution is spawned by a binary integer programming method; then, evolutionary algorithms are designed to achieve a final optimized solution, which minimizes efficiency losses of participants.

Goudarzi et al. [56] provides a solution method for a multi-dimensional capacity allocation problem in a IaaS system, while guaranteeing SLA requirements to customers running multi-tiers applications. An improved solution is obtained starting from an initial configuration based on an upper bound; then, a force-directed search is adopted to increase the total profit. Moreover, a closedform formula for calculating the average response time of a request and a unified framework to manage different levels of SLAs are provided.

Zaman et al. [138] considers an online mechanism for computing resource allocation to VMs subject to limited information. The algorithm evaluates allocation and revenues as the users place requests to the system. Furthermore, the authors prove that their approach is incentive compatible; they also report extensive simulation experiments.

Wang et al. [135] considers capacity allocation subject to two pricing models, a pay-as-you-go offering and periodic auctions. An optimal capacity segmentation strategy is formulated as a Markov decision process, which is then used to maximize revenues. The authors propose also a faster near-optimal algorithm, proven to asymptotically approach the optimal solution, and show a significantly lower complexity with respect to the optimal method.

Roy et al. [18] proposes a model-predictive resource allocation algorithm that auto-scales VMs, with the aim of optimizing the utility of the application over a limited prediction horizon. Empirical results demonstrate that the proposed method satisfies application QoS requirements, while minimizing operational costs.

Dutta et al. [126] develops a resource manager that uses a combination of horizontal and vertical scaling to optimize both resource usage and the reconfiguration cost. Finally, the solution is tested using real production traces.

Zhu et al. [23] builds a VM consolidation algorithm that makes use of an inference model that considers the effect of co-located VMs to predict QoS metrics. In this method, the workload is modeled by means of a Kalman filter, while the resource usage profile is estimated with a Hidden Markov Model. The proposed method is tested against SPECWeb2005.

Hwang et al. [129] also considers the VM consolidation problem by modeling the VM resource demands as a set of correlated random variables. The result is a multi-capacity stochastic bin packing problem, which is solved by means of a simple, scalable yet effective heuristic.

He et al. [128] uses a multivariate probabilistic model to schedule VMs among physical machines in order to improve resource utilization. This approach also considers migration costs, and the multi-dimensional nature of the VM resource requirements (e.g., CPU, memory, and network).

Finally, in [134] a framework that automatically reconfigures the storage system in response to fluctuations in the workload is presented. The framework makes use of a performance model of the system obtained through statistical machine learning. Such model is embedded into an effective Model-Predictive Control algorithm.

#### 4.1.2 Infrastructure-user capacity allocation

From the user perspective, capacity allocation arises in IaaS and PaaS scenarios where the user is in charge with the control of the number of VMs or application containers running in the system. In this context the user is generally a SaaS provider, which wants to maximize her revenues providing a service that meets a certain QoS. Then, the problem to be addressed is to determine the minimum number of VMs or containers needed to fulfill the target QoS, pursuing the best trade-off between cost and performance.

From the user side, capacity allocation is often implemented through auto-scaling policies. Mao and Humphrey [111] defines an auto-scaling mechanism to guarantee the execution of all jobs within given deadlines. The solution accounts for workload burstinesses and delayed instance acquisition. This approach is compared against other techniques and it shows cost savings from 9.8% to 40.4%. Maggio et al. [110] compares several approaches for decision-making, as part of an autonomic framework that allocates resources to a software application.

Patikirikorala et al. [112] proposes a multi-model control-based framework to deal with the highly nonlinear nature of software systems. An extensible meta-model and a class library with an initial set of five models are developed. Finally, the presented approach is endorsed against fixed and adaptive control schemes by means of a campaign of experiments.

In [108] an optimal resource provisioning algorithm is derived to deal with the uncertainty of resource advancereservation. The algorithm reduces resources under- and over-provisioning by minimizing the total cost for a customer during a certain time horizon. The solution methods are based on the Bender decomposition approach to divide the problem into sub-problems, which can be solved in parallel, and an approximation algorithm to solve problems with a large set of scenarios.

On-demand and reserved resources are considered in the model proposed in [107] to define a bio-inspired selfadapting solution for cloud resource provisioning with the aim of minimizing the number of required virtual machines while meeting SLAs.

A decentralized probabilistic algorithm is also described in [106], which focuses on federated clouds. The proposed solution has the aim to take advantage of a Cloud federation to avoid the dependence on a single provider, while still minimizing the amount of used resources to maintain a good QoS level for customers. The solution provides an effective decentralized algorithm for deploying massively scalable services and it is suitable for all the situations in which a centralized solution is not feasible.

Ali-Eldin et al. [114] aims at dynamic resource provisioning exploiting horizontal elasticity. Two adaptive hybrid controllers, including both reactive and proactive actions, are employed to decide the number of VMs for a cloud service to meet the SLAs. The future demand is predicted by a queueing-network model.

A key-value store is presented in [115] to meet lowlatency Service Level Objectives (SLOs). The proposed middleware achieves high scalability by using replication, providing more predictable response times. An analytical model, based on queueing theory, is presented to describe the relation between the number of replicas and the service level, e.g., the percentage of requests processed according to SLOs.

A capacity allocation problem in presented in [113] that exploits both horizontal and vertical elasticity. An integer linear problem is used to calculate an optimized new configuration able to deal with the current workload. However, reconfiguration is executed if the associated overhead cost calculated on a expected stability duration is lower than a certain minimum benefit defined by a human decision maker.

In [117] two multi-objective customer-driven SLAbased resource provisioning algorithms are proposed. The objectives are the minimization of both resource and penalty costs, as well as minimizing SLA violations. The proposed algorithms consider customer profiles and quality parameters to cope with dynamic workloads and heterogeneous cloud resources.

Finally, a profile-based approach for scalability is described in [109], the authors propose a solution based on the definition of platform-independent profiles, which enable the automation of setup and scaling of application servers in order to achieve a just-in-time scalability of the execution environment, as demonstrated with a case study presented in the paper.

#### 4.2 Load balancing

#### 4.2.1 Infrastructure-provider load balancing

Request load-balancing is an increasingly supported feature of cloud offerings. A load balancer dispatches requests from users to servers according to a load dispatching policy. Policies differ for the decision approach and for the amount of information they use. Research work has focused on policies that are either simple to implement, and thus minimize overheads, or that offer some optimality guarantees, typically proven by analytical models.

The research literature has investigated both centralized and decentralized load balancing mechanisms for providers.

Among centralized approaches, [122] introduces an offline optimization problem for geographical load balancing among data centers, explicitly considering SLAs and dynamic electricity prices. This is complemented with an online algorithm to handle the uncertainty in electricity prices. The proposed algorithm is compared against a greedy heuristic method and it shows significant cost savings (around 20-30%).

A load balancer is presented in [50] to assign VMs among geographically-distributed data centers considering predictions on workload, energy prices, and renewable energy generation capacities. Two complementary methods are proposed: an offline deterministic optimization method to be used at design time and an online VM placement, migration and geographical load balancing algorithm for runtime. The authors studied the behavior of both online and offline algorithms by means of a simulation campaign. The results demonstrate that online version of the algorithm performs 8% worse than the offline one because it deals with incomplete information. On the other hand, the analysis also shows that turning on the geographical load balancing has a strong impact on quality of the solutions (between 27% and 40%) of the online algorithm.

Spicuglia et al. [141] proposes an online load balancing policy that considers the inherent VM heterogeneity found in cloud resources. The load balancer uses the number of outstanding requests and the inter-departure times in each VM to dispatch requests to the VM with the shortest expected response time. The authors demonstrate that their solution is able to improve the variance and percentiles of response times with respect to a built-in policy of the Apache web server.

Decentralized methods are considered in [107], which proposes a self-organizing approach to provide robust and scalable solutions for service deployment, resource provisioning, and load balancing in a cloud infrastructure. The algorithm developed has the additional benefit to leverage Cloud elasticity to allocate and deallocate resources to help services to respect contractual SLAs.

Another example is the cost minimization mechanism for data-intensive service provisioning proposed in [121]. Such mechanism uses biological evolution concepts to manage data application services and to produce optimal composition and load balancing solutions. A multiobjective genetic algorithm is described in detail but a systematic experimental campaign is planned as future work.

## 4.2.2 Infrastructure-user load balancing

In the studies considered in the previous section, the load balancer is installed and managed transparently by the cloud provider. In some cases, the user can decide to install its own load balancer for a cloud application. This may be helpful, for instance, to jointly tackle capacity allocation and load balancing.

For example, [47] considers a joint optimization problem on multiple IaaS service centers. A non-linear model for the capacity allocation and load redirection of multiple request classes is proposed and solved by decomposition. A comparison against a set of heuristics from the literature and an oracle with perfect knowledge about the future load shows that the proposed algorithm overcomes the heuristic approaches, without penalizing SLAs and it is able to produce results that are close to the global optimum. Anselmi and Casale [120] provides a simple heuristic for user-side load-balancing under connection pooling that is validated against an IaaS cloud dataset. The main result is that the presented approach is able to provide tight guarantees on the optimality gap and experimental results show that it is at the same time accurate and fast.

Hybrid clouds are considered in [116]. The authors formulate an optimization problem faced by a cloud procurement endpoint (a module responsible for provisioning resources from public cloud providers), where heavy workloads are tackled by relying on public clouds. They present a linear integer program to minimize the resource cost, and evaluate how the solution scales with the different problem parameters.

In [123] a structured peer-to-peer network, based on distributed hash tables, is proposed to support service discovery, self-management, and load-balancing of cloud applications. The effectiveness of the peer-to-peer approach is demonstrated through a set of experiments executed on Amazon EC2.

Finally, [105] proposes an adaptive approach for component replication of cloud applications, aiming at finding a cost-effective placement and load balancing. This is a distributed method based on an economic multiagent model that achieves high application availability guaranteeing at the same time service availability under failures.

#### 4.3 Admission control

#### 4.3.1 Infrastructure-provider admission control

Admission control is an overload protection mechanism that rejects requests under peak workload conditions to prevent QoS degradation. A lot of work has been done in the last decade for optimal admission control in web servers and multi-tier applications. The basic idea is to predict the value of a specific QoS metric and if such value grows above a certain threshold, the admission controller rejects all new sessions favoring the service of requests from already admitted sessions.

In cloud computing, several works on admission control have emerged in IaaS. Khazaei et al. [130] develops an analytical model for resource provisioning, virtual machine deployment, and pool management. This model predicts service delay, task rejection probability, and steady-state distribution of server pools.

The availability of resources and admission control is also discussed in [131]. The work uses a probabilistic approach to find an optimized allocation of services on virtualized physical resources. The main requirement of this system is the horizontal elasticity. In fact, the probability of requesting more resources for a service is at the basis of the formulated optimization model, that constitutes a probabilistic admission control test.

Almeida et al. [49] proposes a joint admission control and capacity allocation algorithm for virtualized IaaS systems minimizing the data center energy costs and the penalty incurred for request rejections and SLA violations. SLAs are expressed in terms of the tail distribution of application response times.

Agostinho et al. [124] optimizes the allocation and scheduling of VMs in federated clouds using a genetic algorithm. The solution is composed by two parts: First, servers selection in a data-center is performed by using a search based bio-inspired technique; then, data centers are selected within the cloud federation by using a shortest path algorithm, according to the available bandwidth of links connecting the domains. The aim of the paper is to exploit resources in domains with low allocation costs and, at the same time, achieve better network performance among cloud nodes.

Ellens et al. [52] allows service providers to reserve a certain amount of resources exclusively for some customers, according to SLAs. The proposed framework helps to stipulate a realistic SLA with customers and supports dynamic load shedding and capacity provisioning by considering a queueing model with multiple priority classes. The main performance metric being optimized is the rejection probability, which has to guarantee the value stipulated in the SLA.

The work in [139] proposes an admission control protocol to prevent over-utilization of system resources, classifying applications based on resource quality requirements. It uses an open multi-class queueing network to support a QoS-aware admission control on heterogeneous resources to increase system throughput.

In order to control overload in Database-as-a-Service (DaaS) environments, [140] proposes a profit-aware admission control policy. It first uses nonlinear regression to predict the probability for a query to meet its requirement, and then decides whether the query should be admitted to the database system or not.

## 4.3.2 Infrastructure-user admission control

From the cloud-user perspective, the admission control mechanism is used as an extreme overload mechanism, helpful when additional resources are obtained with some significant delay. For example, during a cloud burst (i.e., when part of the application traffic is redirected from a private to a public data center to cope with a traffic intensity that surpasses the capacity of the private infrastructure), if the public cloud resources are not provided timely, one can decide to drop new incoming request to preserve the QoS for users already in the system (or at least part of them, e.g., *gold customers*), avoiding application performance degradation.

Three different admission control and scheduling algorithms are proposed in [119] to effectively exploiting public cloud resources. The paper takes the perspective of a SaaS provider with the aim of maximizing the profit by minimizing cost and improving customer satisfaction levels.

Leitner et al. [118] introduces a client-side admission control method to schedule requests among VMs, looking at minimizing the cost of application, SLA violations and IaaS resources.

## 5 Discussion and conclusion

In recent years, cloud computing has matured from an early-stage solution to a mainstream operational model for enterprise applications. However, the diversity of technologies used in cloud systems makes it difficult to analyze their QoS and, from the provider perspective, to offer service-level guarantees. We have surveyed current approaches in workload and system modeling and early applications to cloud QoS management.

From this survey, a number of insights arise on the current state of the art:

- The number of works that apply white-box system modeling techniques is quite limited in QoS management, albeit popular in the software performance engineering community. This effectively creates a divide between the knowledge that can be made available for an application by its designers and the techniques used to manage it. A research question is whether the availability of more detailed information about application internals can provide significant advantages in QoS management. Indeed, a trade-off exists between available information, QoS model complexity, computational cost of decision-making, and accuracy of predictions. This trade-off requires further investigation by the research community.
- Gray-box models that emphasize resource consumption modeling are currently prevalent in QoS management studies. However, description of performance is often quite basic and associated with mean resource

requirements of the applications. However, the cloud measurement studies in Section 2.1 have identified performance variability as a major issue in today's offerings, calling for more comprehensive models that can describe also the variability in CPU requirements, in addition to mean requirements. Such extension has been explored in black-box system models (e.g., QoS in web services), but it is far less understood in white-box and gray-box modeling.

- Quite surprisingly, we have found a limited amount of work specific to workload analysis and inference techniques in the cloud. Most of the techniques used for traffic forecasting, resource consumption estimation, and anomaly detection have received little or no validation in a cloud environment. As such, it remains to establish the robustness of current techniques to noisy measurements typical of multi-tenant cloud environments.
- Another observation arising from our survey is that the literature is rich in works focusing on IaaS systems, often deployed on Amazon EC2, at present the market leader in this segment.
- If we consider the resource management mechanisms for applications QoS enforcement provided by public clouds, they are quite simplistic if compared to current research proposals. Indeed, such mechanisms are mainly *reactive* and are triggered by thresholds violations (related to response times, as in Google App Engine, or CPU utilization or other low level infrastructure metrics, as in Amazon EC2.) Vice versa, integrating workload characterisation, system models and resource management solutions, *pro-active* systems, may help to prevent QoS degradation. The development of research prototypes that are transferable in commercial solutions seems to remain an open point.
- Finally, in cloud systems an important role is played by resource pricing models. There is a growing interest towards understanding better cloud spot markets, where bidding strategies are developed for procuring computing resources. Approaches are currently being proposed to automate dynamic pricing and cloud resources selection. We expect that, in upcoming years, these models will play a bigger role than today in capacity allocation frameworks.

Summarizing, this survey shows that the literature has already a significant number of works in cloud QoS management, but their focus leaves open several research opportunities in the areas discussed above.

#### Endnote

<sup>a</sup> Throughout this paper, we mainly focus on QoS aspects pertaining to performance, reliability and

#### **Competing interests**

The authors declare that they have no competing interests.

#### Authors' contributions

All authors read and approved the final manuscript.

#### Acknowledgment

The research reported in this article is partially supported by the European Commission grant FP7-ICT-2011-8-318484 (MODAClouds, www.modaclouds.eu).

#### Author details

<sup>1</sup> Dipartimento di Elettronica, Informazione, e Bioingegneria Politecnico di Milano, Piazza Leonardo da Vinci 32, Milan 20133, Italy. <sup>2</sup> Department of Computing, Imperial College London, 180 Queens Gate, London SW7 2AZ, UK.

#### Received: 5 February 2014 Accepted: 12 August 2014 Published online: 25 September 2014

#### References

- Armbrust M, Fox A, Griffith R, Joseph AD, Katz R, Konwinski A, Lee G, Patterson D, Rabkin A, Stoica I, Zaharia M (2010) A view of cloud computing. Commun ACM 53(4):50–58
- Zhang Q, Cheng L, Boutaba R (2010) Cloud computing: state-of-the-art and research challenges. J Internet Serv Appl 1(1):7–18
- Ardagna D, Panicucci B, Trubian M, Zhang L (2012) Energy-aware autonomic resource allocation in multitier virtualized environments. IEEE Trans Serv Comput 5(1):2–19
- Petcu D, 0 Macariu G, Panica S, Craciun C (2013) Portable cloud applications - from theory to practice. Future Generation Comput Syst 29(6):1417–1430
- Farley B, Juels A, Varadarajan V, Ristenpart T, Bowers KD, Swift MM (2012) More for your money: Exploiting performance heterogeneity in public clouds. In: Proceedings of the 2012 Third ACM Symposium on Cloud Computing, SoCC '12, San Jose, CA, USA, pp 1–14
- Ou Z, Zhuang H, Nurminen JK, Ylä-Jääski A, Hui P (2012) Exploiting hardware heterogeneity within the same instance type of amazon ec2. In: Proceedings of the 4th USENIX Conference on Hot Topics in Cloud Ccomputing, HotCloud'12, Boston, MA, USA, pp 4–4
- Schad J, Dittrich J, Quiané-Ruiz J-A (2010) Runtime measurements in the cloud: Observing, analyzing, and reducing variance. Proc VLDB Endowment 3(1–2):460–471
- Mao M, Humphrey M (2012) A performance study on the VM startup time in the cloud. In: Proceedings of the 2012 IEEE Fifth International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 423–430
- Xu Y, Musgrave Z, Noble B, Bailey M (2013) Bobtail: Avoiding long tails in the cloud. In: Proceedings of the 10th USENIX Conference on Networked Systems Design and Implementation, NSDI '13, Lombard, IL, USA, pp 329–342
- Wang G, Ng TSE (2010) The impact of virtualization on network performance of amazon ec2 data center. In: Proceedings of the 29th Conference on Information Communications, INFOCOM'10, San Diego, CA, USA, pp 1163–1171
- 11. Hill Z, Li J, Mao M, Ruiz-Alvarez A, Humphrey M (2011) Early observations on the performance of Windows Azure. Sci Program 19(2–3):121–132
- Li Z, O'Brien L, Ranjan R, Zhang M (2013) Early observations on performance of Google compute engine for scientific computing. In: Proceedings of the 2013 IEEE 5th International Conference on Cloud Computing Technology and Science, volume 1 of CloudCom 2013, Bristol, United Kingdom, pp 1–8
- Drago I, Mellia M, Munafo MM, Sperotto A, Sadre R, Pras A (2012) Inside Dropbox: understanding personal cloud storage services. In: Proceedings of the 2012 ACM Conference on Internet Measurement Conference, IMC '12, Boston, MA, USA, pp 481–494

- Kossmann D, Kraska T, Loesing S (2010) An evaluation of alternative architectures for transaction processing in the cloud. In: Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data, SIGMOD '10, Indianapolis, IN, USA, pp 579–590
- Wada H, Fekete A, Zhao L, Lee K, Liu A (2011) Data consistency properties and the trade-offs in commercial cloud storage: the consumers' perspective. In: Proceedings of the 5th Biennial Conference on Innovative Data Systems Research, CIDR 2011, Asilomar, CA, USA, pp 134–143
- Liu S, Huang X, Fu H, Yang G (2013) Understanding data characteristics and access patterns in a cloud storage system. In: Proceedings of the 2013 13th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing, CCGrid 2013, Delft, Nederlands, pp 327–334
- Li A, Yang X, Kandula S, Zhang M (2010) Cloudcmp: comparing public cloud providers. In: Proceedings of the 10th ACM SIGCOMM conference on Internet measurement. ACM, pp 1–14
- Roy N, Dubey A, Gokhale A (2011) Efficient autoscaling in the cloud using predictive models for workload forecasting. In: Proceedings of the 2011 IEEE International Conference on Cloud Computing, CLOUD '11, Washington, DC, USA, pp 500–507
- Gasquet C, Witomski P (1999) Fourier analysis and applications: filtering, numerical computation, wavelets, volume 30 of Texts in applied mathematics. Springer, New York, USA
- 20. Khan A, Yan X, Shu T, Anerousis N (2012) Workload characterization and prediction in the cloud: A multiple time series approach. In: Proceedings of the 2012 IEEE Network Operations and Management Symposium, NOMS 2012, Maui, HI, USA, pp 1287–1294
- Di S, Kondo D, Walfredo C (2012) Host load prediction in a google compute cloud with a Bayesian model. In: Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, SC12, Salt Lake City, Utah, USA, pp 1–11
- Gmach D, Rolia J, Cherkasova L, Kemper A (2007) Workload analysis and demand prediction of enterprise data center applications. In: Proceedings of the 2007 IEEE 10th International Symposium on Workload Characterization, IISWC '07, Boston, MA, USA, pp 171–180
- Zhu Q, Tung T (2012) A performance interference model for managing consolidated workloads in QoS-aware clouds. In: Proceedings of the 2012 IEEE Fifth International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 170–179
- 24. Hoffmann GA, Trivedi KS, Malek M (2007) A best practice guide to resource forecasting for computing systems. IEEE Trans Reliability 56(4):615–628
- Anandkumar A, Bisdikian C, Agrawal D (2008) Tracking in a spaghetti bowl: Monitoring transactions using footprints. In: Proceedings of the 2008 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems. ACM Press, Annapolis, Maryland, USA, pp 133–144
- 26. Menascé D, Almeida V, Dowdy L (1994) Capacity planning and performance modeling: from mainframes to client-server systems. Prentice-Hall, Inc. NJ, USA
- Rolia J, Vetland V (1995) Parameter estimation for performance models of distributed application systems. In: In Proc. of CASCON. IBM Press, Toronto, Ontario, Canada, p 54
- Rolia J, Vetland V (1998) Correlating resource demand information with ARM data for application services. In: Proceedings of the 1st international workshop on Software and performance. ACM, Santa Fe, New Mexico, USA, pp 219–230
- 29. Liu Y, Gorton I, Fekete A (2005) Design-level performance prediction of component-based applications. IEEE Trans Softw Eng 31(11):928–941
- Sutton CA, Jordan MI (2008) Probabilistic inference in queueing networks. In: Proceedings of the 3rd conference on Tackling computer systems problems with machine learning techniques. USENIX Association, Berkeley, CA, US, p 6
- Sutton CA, Jordan MI (2010) Inference and learning in networks of queues. In: International Conference on Artificial Intelligence and Statistics, Sardinia, Italy, pp 796–803
- Zhang Q, Cherkasova L, Smirni E (2007) A regression-based analytic model for dynamic resource provisioning of multi-tier applications. In: Proc. of the 4th ICAC Conference, Jacksonville, Florida, USA, pp 27–27
- Liu Z, Wynter L, Xia C, Zhang F (2006) Parameter inference of queueing models for it systems using end-to-end measurements. Perform Eval 63(1):36–60

- Casale G, Cremonesi P, Turrin R (2008) Robust workload estimation in queueing network performance models. In Proc. of Euromicro PDP:183–187
- 35. Kalbasi A, Krishnamurthy D, Rolia J, Dawson S (2012) DEC: Service demand estimation with confidence. IEEE Trans Softw Eng 38(3):561–578
- Kalbasi A, Krishnamurthy D, Rolia J, Richter M (2011) MODE: Mix driven on-line resource demand estimation. In: Proceedings of the 7th International Conference on Network and Services Management. International Federation for Information Processing, pp 1–9
- Pacifici G, Segmuller W, Spreitzer M, Tantawi A (2008) CPU demand for web serving: Measurement analysis and dynamic estimation. Perform Eval 65(6):531–553
- Cremonesi P, Sansottera A (2012) Indirect estimation of service demands in the presence of structural changes. In: Proceedings of Quantitative Evaluation of Systems (QEST). IEEE, London, UK, pp 249–259
- Cremonesi P, Dhyani K, Sansottera A (2010) Service time estimation with a refinement enhanced hybrid clustering algorithm. In: Analytical and Stochastic Modeling Techniques and Applications. Springer, pp 291–305
- Sharma AB, Bhagwan R, Choudhury M, Golubchik L, Govindan R, Voelker GM (2008) Automatic request categorization in internet services. ACM SIGMETRICS Perform Eval Rev 36(2):16–25
- Wu X, Woodside M (2008) A calibration framework for capturing and calibrating software performance models. In: Computer Performance Engineering. Springer, pp 32–47
- 42. Zheng T, Woodside CM, Litoiu M (2008) Performance model estimation and tracking using optimal filters. IEEE Trans Softw Eng 34(3):391–406
- Desnoyers P, Wood T, Shenoy PJ, Singh R, Patil S, Vin HM (2012) Modellus: Automated modeling of complex internet data center applications. TWEB 6(2):8
- Longo F, Ghosh R, Naik VK, Trivedi KS (2011) A scalable availability model for Infrastructure-as-a-Service cloud. In: Proceedings of the 2011 IEEE/IFIP 41st International Conference on Dependable Systems Networks, DSN 2011, Hong Kong, China, pp 335–346
- Calinescu R, Ghezzi C, Kwiatkowska MZ, Mirandola R (2012) Self-adaptive software needs quantitative verification at runtime. Commun ACM 55(9):69–77
- 46. Chung M-Y, Ciardo G, Donatelli S, He N, Plateau B, Stewart W, Sulaiman E, Yu J (2004) Comparison of structural formalisms for modeling large markov models. In: Parallel and Distributed Processing Symposium, 2004 Proceedings. 18th International, Santa Fe, New Mexico, USA, p 196
- Ardagna D, Casolari S, Colajanni M, Panicucci B (2012) Dual time-scale distributed capacity allocation and load redirect algorithms for cloud systems. J Parallel Distributed Comput 72(6):796–808
- Xiong P, Wang Z, Malkowski S, Wang Q, Jayasinghe D, Pu C (2011) Economical and robust provisioning of n-tier cloud workloads: A multi-level control approach. In: Proceedings of the 31st IEEE International Conference on Distributed Computing Systems (ICDCS), Minneapolis, Minnesota, USA, pp 571–580
- Almeida J, Almeida V, Ardagna D, Cunha I, Francalanci C, Trubian M (2010) Joint admission control and resource allocation in virtualized servers. J Parallel Distributed Comput 70(4):344–362
- Goudarzi H, Pedram M (2013) Geographical load balancing for online service applications in distributed datacenters. In: Proceedings of the 2013 IEEE Sixth International Conference on Cloud Computing, CLOUD '13, Santa Clara, CA, USA, pp 351–358
- Zhang Q, Zhu Q, Zhani MF, Boutaba R (2012) Dynamic service placement in geographically distributed clouds. In: Proceedings of the 2012 IEEE 32Nd International Conference on Distributed Computing Systems, ICDCS '12, Macau, China, pp 526–535
- Ellens W, Zivkovic M, Akkerboom J, Litjens R, van den Berg H (2012) Performance of cloud computing centers with multiple priority classes. In: Proceedings of the 2012 IEEE 5th International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 245–252
- Kusic D, Kandasamy N (2006) Risk-aware limited lookahead control for dynamic resource provisioning in enterprise computing systems. In: Proceedings of the 2006 IEEE International Conference on Autonomic Computing, ICAC '06, Dublin, Ireland, pp 74–83
- Kusic D, Kephart JO, Hanson JE, Kandasamy N, Jiang G (2009) Power and performance management of virtualized computing environments via lookahead control. Cluster Comput 12(1):1–15

- Addis B, Ardagna D, Panicucci B, Squillante MS, Zhang L (2013) A hierarchical approach for the resource management of very large cloud platforms. IEEE Trans Dependable Secure Comput 10(5):253–272
- Goudarzi H, Pedram M (2011) Multi-dimensional sla-based resource allocation for multi-tier cloud computing systems. In: Proceedings of the 2011 IEEE 4th International Conference on Cloud Computing, CLOUD '11, Washington, DC, USA, pp 324–331
- 57. Becker 5, Koziolek H, Reussner R (2009) The Palladio component model for model-driven performance prediction. J Syst Softw 82(1):3–22
- Franks G, Al-Omari T, Woodside CM, Das O, Derisavi S (2009) Enhanced modeling and solution of layered queueing networks. IEEE Trans Softw Eng 35(2):148–161
- Omari T, Franks G, Woodside M, Pan A (2007) Efficient performance models for layered server systems with replicated servers and parallel behaviour. J Syst Softw 80(4):510–527
- 60. Tribastone M (2013) A fluid model for layered queueing networks. IEEE Trans Softw Eng 39(6):744–756
- Pérez JF, Casale G (2013) Assessing sla compliance from palladio component models. In: Proceedings of the 2013 15th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing, SYNASC '13, Timisoara, Romania, pp 409–416
- Faisal A, Petriu D, Woodside M (2013) Network latency impact on performance of software deployed across multiple clouds. In: Proceedings of the 2013 Conference of the Center for Advanced Studies on Collaborative Research, CASCON '13, Ontario, Canada, pp 216–229
- 63. Jung G, Joshi KR, Hiltunen MA, Schlichting RD, Pu C (2008) Generating adaptation policies for multi-tier applications in consolidated server environments. In: Autonomic Computing, 2008 ICAC'08. International Conference on. IEEE, Chicago, IL, USA, pp 23–32
- 64. Bacigalupo D, van Hemert J, Chen X, Usmani A, Chester A, He L Dillenberger D, Wills G, Gilbert L, Jarvis S (2011) Managing dynamic enterprise and urgent workloads on clouds using layered queuing and historical performance models. Simul Model Prac Theory 19:1479–1495
- 65. Thereska E, Ganger GR (2008) IRONmodel: Robust performance models in the wild. ACM SIGMETRICS Perform Eval Rev 36(1):253–264
- Singh R, Sharma U, Cecchet E, Shenoy P (2010) Autonomic mix-aware provisioning for non-stationary data center workloads. In: Proceedings of the 7th ACM international conference on Autonomic computing, Washington, DC, USA, pp 21–30
- Brogi A, Corfini S, Iardella S (2009) From OWL-S descriptions to Petri nets. In: Service-Oriented Computing - ICSOC 2007 Workshops, volume 4907 of Lecture Notes in Computer Science. Springer-Verlag, Vienna, Austria, pp 427–438
- Callou G, Maciel P, Tutsch D, Araujo J (2011) Models for dependability and sustainability analysis of data center cooling architectures. In: Proceedings of the 2011 Symposium on Theory of Modeling & amp; Simulation: DEVS Integrative M& amp; S Symposium, TMS-DEVS '11, Boston, MA, USA, pp 274–281
- 69. Wei B, Lin C, Kong X (2011) Dependability modeling and analysis for the virtual data center of cloud computing. In: Proceedings of the 2011 IEEE 13th International Conference on High Performance Computing and Communications, HPCC 2011, Bamff, Canada, pp 784–789
- Figueiredo J, Maciel P, Callou G, Tavares E, Sousa E, Silva B (2011) Estimating reliability importance and total cost of acquisition for data center power infrastructures. In: Proceedings of the 2011 IEEE International Conference on Systems, Man, and Cybernetics, SMC 2011, Anchorage, AK, USA, pp 421–426
- Dantas J, Matos R, Araujo J, Maciel P (2012) An availability model for eucalyptus platform: An analysis of warm-standy replication mechanism. In: Proceedings of the 2012 IEEE International Conference on Systems, Man, and Cybernetics, SMC 2012, Seoul, Korea, pp 1664–1669
- Melo M, Maciel P, Araujo J, Matos R, Araujo C (2013) Availability study on cloud computing environments: Live migration as a rejuvenation mechanism. In: Proceedings of the 2013 IEEE/IFIP 43rd International Conference on Dependable Systems and Networks, DSN 2013, Hong Kong, China, pp 1–6
- Ford B (2012) Icebergs in the clouds: The other risks of cloud computing. In: Proceedings of the 2012 4th USENIX Conference on Hot Topics in Cloud Computing, HotCloud'12, Boston, MA, USA, pp 2–2
- Jhawar R, Piuri V (2012) Fault tolerance management in laaS clouds. In: Proceedings of 2012 IEEE First AESS European Conference on Satellite Telecommunications, ESTEL 2012, Rome, Italy, pp 1–6

- Kiran M, Jiang M, Armstrong DJ, Djemame K (2011) Towards a service lifecycle based methodology for risk assessment in cloud computing. In: Proceedings of the 2011 IEEE Ninth International Conference on Dependable, Autonomic and Secure Computing, DASC 2011, Sydney, NSW, Australia, pp 449–456
- Ardagna D, Pernici B (2007) Adaptive service composition in flexible processes. IEEE Trans Softw Eng 33(6):369–384
- Ben Mabrouk N, Beauche S, Kuznetsova E, Georgantas N, Issarny V (2009) QoS-aware service composition in dynamic service oriented environments. In: Proceedings of the 10th ACM/IFIP/USENIX International Conference on Middleware, Middleware '09, Urbanna, II, USA, pp 1–20
- Alrifai M, Skoutas D, Risse T (2010) Selecting skyline services for QoS-based web service composition. In: Proceedings of the 19th International Conference on World Wide Web, WWW '10, Raleigh, NC, USA, pp 11–20
- El Haddad J, Manouvrier M, Rukoz M (2010) TQoS: Transactional and QoS-aware selection algorithm for automatic web service composition. IEEE Trans Serv Comput 3(1):73–85
- Cardellini V, Casalicchio E, Grassi V, Mirandola R (2006) A framework for optimal service selection in broker-based architectures with multiple QoS classes. In: Proceedings of the 2006 IEEE Services Computing Workshops, SCW '06, Chicago, IL, USA, pp 105–112
- Jiang D, Pierre G, Chi C-H (2010) Autonomous resource provisioning for multi-service web applications. In: Proceedings of the 2010 19th International Conference on World Wide Web, WWW '10, Raleigh, NC, USA, pp 471–480
- Yu T, Zhang Y, Lin K-J (2007) Efficient algorithms for web services selection with end-to-end QoS constraints. ACM Trans Web 1(1):6
- Stein S, Payne TR, Jennings NR (2009) Flexible provisioning of web service workflows. ACM Trans Internet Technol 9(1):1–45
- Comuzzi M, Pernici B (2009) A framework for qos-based web service contracting. ACM Trans Web 3(3):1–52
- Schuller D, Polyvyanyy A, García-Bañuelos L, Schulte S (2011) Optimization of complex qos-aware service compositions. In: Proceedings of the 2011 9th International Conference on Service-Oriented Computing, ICSOC'11, Paphos, Cyprus, pp 452–466
- Clark A, Gilmore S, Tribastone M (2009) Quantitative analysis of web services using srmc. In: Formal Methods for Web Services, volume 5569 of Lecture Notes in Computer Science. Springer, Berlin Heidelberg, pp 296–339
- Reinecke P, Wolter K (2008) Phase-type approximations for message transmission times in web services reliable messaging. In: Proceedings of the 2008 SPEC International Workshop on Performance Evaluation: Metrics, Models and Benchmarks, SIPEW '08 Boston, MA, USA, Darmstadt, Germany, pp 191–207
- Haddad S, Mokdad L, Youcef S (2010) Response time of BPEL4WS constructors. In: Proceedings of the 2010 IEEE Symposium on Computers and Communications, ISCC'10, Riccione, Italy, pp 695–700
- Menascé DA, Casalicchio E, Dubey VK (2010) On optimal service selection in service oriented architectures. Perform Eval 67(8):659–675
- Huang J, Liu Y, Duan Q (2012) Service provisioning in virtualization-based cloud computing: Modeling and optimization. In: Proceedings of 2012 IEEE Global Communications Conference, GLOBECOM 2012, Anaheim, CA, USA, pp 1710–1715
- Klein A, Ishikawa F, Honiden S (2012) Towards network-aware service composition in the cloud. In: Proceedings of the 21st International Conference on World Wide Web, WWW '12, Lyon, France, pp 959–968
- Ye Z, Bouguettaya A, Zhou X (2012) QoS-aware cloud service composition based on economic models. In: Proceedings of the 10th International Conference on Service-Oriented Computing, ICSOC'12, Shanghai, China, pp 111–126
- Calheiros RN, Ranjan R, Beloglazov A, De Rose CAF, Buyya R (2011) Cloudsim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms. Software–Pract Exp 41(1):23–50
- 94. Wickremasinghe B, Calheiros RN, Buyya R (2010) CloudAnalyst: A cloudsim-based visual modeller for analysing cloud computing environments and applications. In: Proceedings of the 2010 24th IEEE International Conference on Advanced Information Networking and Applications, AINA 2010, Perth, Australia, pp 446–452
- 95. Calheiros RN, Netto MAS, De Rose CAF, Buyya R (2013) Emusim: an integrated emulation and simulation environment for modeling,

evaluation, validation of performance of cloud computing applications. Software–Pract Exp 43:595–612

- 96. Calheiros RN, Buyya R, De Rose CAF (2010) Building an automated and self-configurable emulation testbed for grid applications. Software–Pract Exp 40:405–429
- Kliazovich D, Bouvry P, Khan SU (2010) GreenCloud: A packet-level simulator of energy-aware cloud computing data centers. In: Proceedings of the 2010 IEEE Global Telecommunications Conference, GLOBECOM 2010, Miami, FL, USA, pp 1–5
- 98. The Network Simulator NS. (http://www.isi.edu/nsnam/ns/)
- Keller G, Tighe M, Lutfiyya H, Bauer M (2012) DCSim: A data centre simulation tool. In: Proceedings of 2012 8th international conference on Network and service management, and 2012 workshop on systems virtualiztion management, CNSM-SVM 2012, Las Vegas, NV, USA, pp 385–392
- Ostermann S, Plankensteiner K, Prodan R, Fahringer T (2011) Groudsim: An event-based simulation framework for computational grids and clouds. In: Proceedings of the 2010 Conference on Parallel Processing, Euro-Par 2010, Ischia, Italy, pp 305–313
- 101. Wang H, Jing Q, Chen R, He B, Qian Z, Zhou L (2010) Distributed systems meet economics: pricing in the cloud. In: Proceedings of the 2nd USENIX conference on Hot topics in cloud computing, HotCloud'10, Boston, MA, USA, pp 6–6
- 102. Sowmya K, Sundarraj RP (2012) Strategic bidding for cloud resources under dynamic pricing schemes. In: Proceedings of 2012 International Symposium on Cloud and Services Computing, ISCOS 2012, Mangalore, India, pp 25–30
- Beloglazov A, Buyya R, Lee YC, Zomaya AY (2011) A taxonomy and survey of energy-efficient data centers and cloud computing systems. Adv Comput 82:47–111
- Aleti A, Buhnova B, Grunske L, Koziolek A, Meedeniya I (2013) Software architecture optimization methods: A systematic literature review. IEEE Trans Softw Eng 39(5):658–683
- Bonvin N, Papaioannou T, Aberer K (2010) An economic approach for scalable and highly-available distributed applications. In: Proceedings of the 2010 IEEE 3rd International Conference on Cloud Computing, CLOUD '10, Miami, FL, USA, pp 498–505
- Calcavecchia NM, Caprarescu BA, Di Nitto E, Dubois DJ, Petcu D (2012) DEPAS: A decentralized probabilistic algorithm for auto-scaling. Computing 94(8–10)
- 107. Caprarescu BA, Calcavecchia NM, Di Nitto E, Dubois DJ (2012) Sos cloud: Self-organizing services in the cloud. In: Bio-Inspired Models of Network, Information, and Computing Systems, volume 87 of Lecture Notes of the Institute for Computer Sciences, Social Informatics and Telecommunications Engineering. Springer, Berlin Heidelberg, pp 48–55
- Chaisiri S, Lee B-S, Niyato D (2012) Optimization of resource provisioning cost in cloud computing. IEEE Trans Serv Comput 5(2):164–177
- 109. Jie Y, Jie Q, Ying L (2009) A profile-based approach to just-in-time scalability for cloud applications. In: Proceedings of the 2009 IEEE International Conference on Cloud Computing, CLOUD '09, Bangalore, India, pp 9–16
- Maggio M, Hoffmann H, Santambrogio MD, Agarwal A, Leva A (2011) A comparison of autonomic decision making techniques. Technical Report MIT-CSAIL-TR-2011-019, Massachusetts Institute of Technology. USA, Massachusetts
- 111. Mao M, Humphrey M (2011) Auto-scaling to minimize cost and meet application deadlines in cloud workflows. In: Proceedings of the 2011 International Conference for High Performance Computing, Networking, Storage and Analysis, SC '11, Seattle, WA, USA, pp 1–12
- 112. Patikirikorala T, Colman A, Han J, Wang L (2011) A multi-model framework to implement self-managing control systems for qos management. In: Proceedings of the 6th International Symposium on Software Engineering for Adaptive and Self-Managing Systems, SEAMS '11, Honolulu, HI, USA, pp 218–227
- 113. Sedaghat M, Hernandez-Rodriguez F, Elmroth E (2013) A virtual machine re-packing approach to the horizontal vs. vertical elasticity trade-off for cloud autoscaling. In: Proceedings of the 2013 ACM Cloud and Autonomic Computing Conference, CAC '13, Miami, FL, USA, pp 6:1–6:10

- 114. Ali-Eldin A, Tordsson J, Elmroth E (2012) An adaptive hybrid elasticity controller for cloud infrastructures. In IEEE Network Operations and Management Symposium (NOMS):204–212
- 115. Stewart C, Chakrabarti A, Griffith R (2013) Zoolander: Efficiently meeting very strict, low-latency slos. In Proceedings of the 10th International Conference on Autonomic Computing (ICAC):265–277
- 116. Van den, Bossche R, Vanmechelen K, Broeckhove J (2010) Cost-optimal scheduling in hybrid iaas clouds for deadline constrained workloads. In: Proceedings of the 2010 IEEE 3rd International Conference on Cloud Computing, CLOUD'10, Miami, FL, USA, pp 228–235
- 117. Wu L, Garg SK, Versteeg S, Buyya R (2013) SLA-based resource provisioning for hosted software as a service applications in cloud computing environments. IEEE Trans Serv Comput 99:1
- 118. Leitner P, Hummer W, Satzger B, Inzinger C, Dustdar S (2012) Cost-efficient and application sla-aware client side request scheduling in an infrastructure-as-a-service cloud. In: Proceedings of the 2012 IEEE 5th International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 213–220
- Wu L, Garg SK, Buyya R (2012) SLA-based admission control for a software-as-a-service provider in cloud computing environments. J Comput Syst Sci 78(5):1280–1299
- 120. Anselmi J, Casale G (2013) Heavy-traffic revenue maximization in parallel multiclass queues. Perform Eval 70(10):806–821
- Wang L, Shen J (2012) Towards bio-inspired cost minimisation for data-intensive service provision. In: Proceedings of the 2012 IEEE First International Conference on Services Economics, SE 2012, Honolulu, HI, USA, pp 16–23
- 122. Adnan MA, Sugihara R, Gupta RK (2012) Energy efficient geographical load balancing via dynamic deferral of workload. In: Proceedings of the 2012 IEEE Fifth International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 188–195
- Ranjan R, Zhao L, Wu X, Liu A, Quiroz A, Parashar M (2010) Peer-to-peer cloud provisioning: Service discovery and load-balancing. In: Cloud Computing, Computer Communications and Networks. Springer, London, pp 195–217
- 124. Agostinho L, Feliciano G, Olivi L, Cardozo E, Guimaraes E (2011) A bio-inspired approach to provisioning of virtual resources in federated clouds. In: Proceedings of the 2011 IEEE Ninth International Conference on Dependable, Autonomic and Secure Computing, DASC 2011, Sydney, NSW, Australia, pp 598–604
- 125. Chen W, Qiao X, Wei J, Huang T (2012) A profit-aware virtual machine deployment optimization framework for cloud platform providers. In: Proceedings of the 2012 IEEE Fift International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 17–24
- 126. Dutta S, Gera S, Verma A, Viswanathan B (2012) Smartscale: Automatic application scaling in enterprise clouds. In: Proceedings of the 2012 IEEE 5th International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 221–228
- 127. Hadji M, Zeghlache D (2012) Minimum cost maximum flow algorithm for dynamic resource allocation in clouds. In: Proceedings of the 2012 IEEE 5th International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 876–882
- 128. He S, Guo L, Ghanem M, Guo Y (2012) Improving resource utilisation in the cloud environment using multivariate probabilistic models. In: Proceedings of the 2012 IEEE 5th International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 574–581
- 129. Hwang I, Pedram M (2013) Hierarchical virtual machine consolidation in a cloud computing system. In: Proceedings of the 2013 IEEE Sixth International Conference on Cloud Computing, CLOUD '13, Santa Clara, CA, USA, pp 196–203
- Khazaei H, Misic J, Misic V, Rashwand S (2013) Analysis of a pool management scheme for cloud computing centers. IEEE Trans Parallel Distributed Syst 24(5):849–861
- Konstanteli K, Cucinotta T, Psychas K, Varvarigou T (2012) Admission control for elastic cloud services. In: Proceedings of the 2012 IEEE 5th International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 41–48
- 132. Meng X, Pappas V, Zhang L (2010) Improving the scalability of data center networks with traffic-aware virtual machine placement. In: Proceedings of the 29th Conference on Information Communications, INFOCOM'10, San Diego, CA, USA, pp 1154–1162

- Tang C, Steinder M, Spreitzer M, Pacifici G (2007) A scalable application placement controller for enterprise data centers. In: Proceedings of the 16th International Conference on World Wide Web, WWW '07, Banff, Canada, pp 331–340
- 134. Trushkowsky B, Bodík P, Fox A, Franklin MJ, Jordan MI, Patterson DA (2011) The SCADS director: Scaling a distributed storage system under stringent performance requirements. In: Proceedings of the 9th USENIX Conference on File and Stroage Technologies, FAST'11, San Jose, CA, USA, pp 12–12
- 135. Wang W, Li B, Liang B (2012) Towards optimal capacity segmentation with hybrid cloud pricing. In: Proceedings of the 2012 IEEE 32nd International Conference on Distributed Computing Systems, ICDCS 2012, Macau, China, pp 425–434
- Wei G, Vasilakos AV, Zheng Y, Xiong N (2010) A game-theoretic method of fair resource allocation for cloud computing services. J Supercomput 54(2):252–269
- 137. Xiao Z, Chen Q, Luo H (2014) Automatic scaling of internet applications for cloud computing services. IEEE Trans Comput 63(5):1111–1123
- Zaman S, Grosu D (2012) An online mechanism for dynamic vm provisioning and allocation in clouds. In: Proceedings of the 2012 IEEE 5th International Conference on Cloud Computing, CLOUD '12, Honolulu, HI, USA, pp 253–260
- 139. Delimitrou C, Bambos N, Kozyrakis C (2013) QoS-aware admission control in heterogeneous datacenters. In: Proceedings of the 2013 10th International Conference on Autonomic Computing, ICAC ÃŢ13, San Jose, CA, USA, pp 291–296
- 140. Xiong P, Chi Y, Zhu S, Tatemura J, Pu C, Hacigümüş H (2011) Activesla: A profit-oriented admission control framework for database-as-aservice providers. In: Proceedings of the 2nd ACM Symposium on Cloud Computing, SOCC '11, Cascais, Portugal, pp 1–14
- 141. Spicuglia S, Chen LY, Binder W (2013) Join the best queue: Reducing performance variability in heterogeneous systems. In: Proceedings of the 2013 IEEE Sixth International Conference on Cloud Computing, CLOUD '13, Santa Clara, CA, USA, pp 139–146
- 142. Huang DT, Niyato D, Wang P (2012) Optimal admission control policy for mobile cloud computing hotspot with cloudlet. In: Proceedings of the 2012 IEEE Wireless Communications and Networking Conference, WCNC 2012, Paris, France, pp 3145–3149
- 143. Padala P, Shin KG, Zhu X, Uysal M, Wang Z, Singhal S, Merchant A, Salem K (2007) Adaptive control of virtualized resources in utility computing environments. In: Proceedings of the 2Nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007, EuroSys'07, Lisbon, Portugal, pp 289–302
- 144. Kephart J, Chan H, Das R, Levine D, Tesauro G, Rawson F, Lefurgy C (2007) Coordinating multiple autonomic managers to achieve specified power-performance tradeoffs. In: Proceedings of the Fourth International Conference on Autonomic Computing, ICAC '07, Jacksonville, FL, USA, pp 24–24

#### doi:10.1186/s13174-014-0011-3

Cite this article as: Ardagna *et al.*: Quality-of-service in cloud computing: modeling techniques and their applications. *Journal of Internet Services and Applications* 2014 5:11.

## Submit your manuscript to a SpringerOpen<sup>®</sup>

- journal and benefit from:
- Convenient online submission
- Rigorous peer review
- Immediate publication on acceptance
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at > springeropen.com