

The Challenge of Cloud Control

Maria Kihl

Dept. of Electrical and Information Technology, Lund University, Sweden

Erik Elmroth, Johan Tordsson

Dept. of Computer Science, Umeå University, Sweden

Karl Erik Årzén, Anders Robertsson

Dept of Automatic Control, Lund University, Sweden

Abstract

Today's cloud data center infrastructures are not even near being able to cope with the enormous and rapidly varying capacity demands that will be reality in a near future. So far, very little is understood about how to transform today's data centers (being large, power-hungry facilities, and operated through heroic efforts by numerous administrators) into a self-managed, dynamic, and dependable infrastructure, constantly delivering expected QoS with reasonable operation costs and acceptable carbon footprint for large-scale services with sometimes dramatic variations in capacity demands. In this paper, we discuss some of the major challenges for resource-optimized cloud data center. We propose a new research area called Cloud Control, which is a control theoretic approach to a range of cloud management problems, aiming to transform today's static and energy consuming cloud data centers into self-managed, dynamic, and dependable infrastructures, constantly delivering expected quality of service with acceptable operation costs and carbon footprint for large-scale services with varying capacity demands.

1. Introduction

During recent years, several dramatic events have shown that the capacity requirements of online services can increase enormously within minutes. Extreme examples are the CNN web-site during early reporting of the 9/11 terror attack (load doubled every 7th minute with peak 2000% over normal) and the 2500% load increase at Al-Jazeera's web-site during the fourth day of the Egyptian revolution, and the massive peak just minutes after president Mubarak resigned after 18 days.

With cloud data centers currently growing to tens and hundreds of thousands of interconnected servers and the increasing complexity of interconnected applications, system management is growing in complexity to a scale necessitating new behavioral abstractions and models for autonomic computing [1]. Due to the nonlinearity of emergent local behavior, it is intrinsically challenging to understand the correlations between local and global behavior and the effects of local and global management actions. Hence, new behavioral and managerial abstractions are needed, together with methods for distributed or hierarchical control.

By combining virtualization, distributed systems, optimization techniques, accurate performance models and autonomic computing, our vision is that today's cloud technology can be turned into a self-managed optimized

infrastructure providing unprecedented gains in terms of cost-efficiency, flexibility, robustness, eco-efficiency and sustainability [2]. In order to facilitate management, virtual machines (VMs) are typically used as a basic building block, allowing management tools to migrate, start, stop, and hibernate instances. Depending on the cloud deployment scenario, the physical infrastructure that hosts the VMs can be internal, external or a combination of both. In either case, it should appear to the cloud user as a single system always delivering sufficient capacity.

2. Challenges

In this paper, we address a range of important inter-related and algorithmically challenging management and control problems focused on cloud data centers, as illustrated in Figure 1.

The first challenge is *Performance models*, which are of paramount importance for the design and development of robust control systems. Notably, cloud data centers are much larger in scale and with a larger variety of workload dynamics than previous telecom and Internet systems. Previous results on server systems have shown that the dynamics may be captured using a black-box approach and rather simple queue models [3] or using flow models [4]. However, for clouds, only a few per-

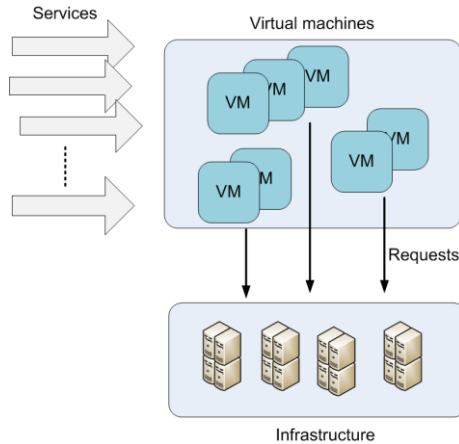


Figure 1 An illustration of a cloud data center.

formance modeling attempts have been presented, see for example [5].

The second challenge is *Service admission control*, which is in essence long term capacity planning where the control system decides to admit a service or not based on anticipated profit. For clouds, only little work has been performed [6].

Further, *Elasticity controllers* should allocate enough resources to a running application to provide acceptable QoS while avoiding costly over-provisioning. Few current solutions actually go beyond simple reactive (non-predictive) threshold-based allocation adjustment [7]. These solutions fall into three broad categories; solutions based on machine learning algorithms, solutions based on control theory and solutions based on statistical workload analysis.

Also, there are a number of *VM placement problems* to be solved to determine where to deploy new VMs and which VMs to shut down [8]. Given the dynamic nature of clouds, with significant changes over time both in demand and supply, VM placement decisions need to be renewed frequently. Recent work on re-placement includes performance modeling of VM migration and performance steering by restricting the number of VMs that can be migrated concurrently. Other dynamic approaches include use of stochastic integer programming to handle uncertainty and genetic algorithms to reduce re-migrations upon load fluctuations.

Another challenge is that large data centers providing cloud services come not only with the cost of buying the equipment but also with a substantial increase in energy costs, which means that *energy optimization* will be crucial [9]. Also, new cyber-physical system categories are emerging where energy interactions between the computing and physical system components are non-negligible [10]. It is important to understand, control, and optimize the energy/temperature interactions at

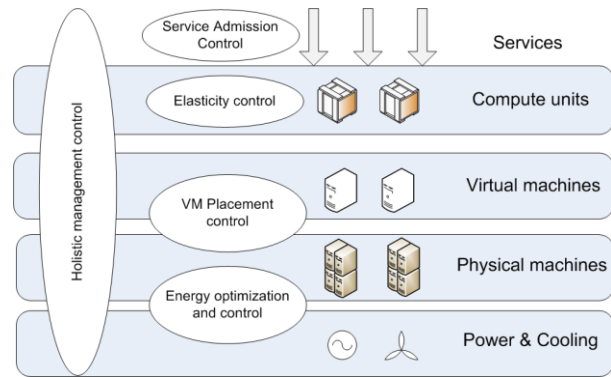


Figure 2 The proposed research area of Cloud Control.

machine room scale (i.e., not only amount of resources but also physical distribution), since these issues are important prerequisites to reducing the energy cost and carbon footprint of cloud systems.

Finally, orchestrating the work of all these different controllers (“*holistic management*”) is a difficult challenge that has not received much attention [11]. Given that workloads and performance models are well understood and that “control knobs” are available in terms of sophisticated methods for performing elasticity control, admission control, VM placement, and power management, there is still a challenge to optimize the whole system behavior. Since the “knobs” may interact in many ways and different performance policies may interfere with one another, uncoordinated interference of multi-knobs can lead to instability or poor performance, calling for a holistic view of resource management [12].

Cloud control

The challenges for resource optimized cloud data centers are immense, and we therefore propose a new research area that we call Cloud Control. The area will address the whole range of closely related challenges for the efficient autonomic management of resources in cloud data centers. Focus is on fundamental modeling, control methods, and algorithms for performing management tasks. The topics are individually important but also fit into the same bigger picture, as illustrated in Figure 2. Each challenge described before can be incorporated in a framework for cloud datacenter resource optimization and control, and some examples will be presented below.

First, the service admission controller decides whether to accept the elastic service that over time may have largely varying capacity demands. A cloud provider must determine the optimal number of services to admit in order to maximize its utility (revenue, utilization, etc.) without endangering Service Level Agreements (SLAs) of already provisioned services. Each admitted

service adds a stochastic load to the cloud infrastructure. Admission control schemes are thus needed to assess the long term effects of new services depending on their estimated workloads.

Second, the objective of the elasticity control system is to allocate enough so called Compute Units (CUs) in order to comply with certain control objectives, targeting the performance expectations of the system. We proposed to define the term CU in order to have a metric for capacity needs, which is independent of the definition of VMs. CUs will later be mapped to suitable VMs. If it is assumed that the system state and workload can be estimated with some accuracy, the controller design basically becomes a stochastic optimal control problem.

Third, the VM placement controller both perform the mapping from CUs to an appropriate set of VMs to run locally or in remote data centers, and determines the packing of VMs on physical machines. We consider the mappings from CU to VM and from VMs to datacenters before performing the assignment of VMs to PMs within the datacenter. These capacity allocation (assignment) problems that can be formulated, e.g., as constrained linear problems for optimizing cost (monetary or other) subject to a capacity and possibly other constraints, such that a specific load balancing (dispersion) objective or requirements on affinity and anti-affinity (e.g., in the same datacenter (or host), not in the same datacenter).

Fourth, the data center energy optimizer models and optimizes the energy-temperature interactions at machine room scale for VM placement decisions. The energy optimization will require a cyber-physical perspective.

Finally, the holistic management system monitors the complete system behavior and optimizes the management tools' concerted operation with respect to the overall data center management objectives.

Conclusions

Cloud computing has received much attention the last years, and is envisioned to be used for a wide range of services. However, the large amount of data centers required for these cloud services will have immense challenges of resource management and control. In this paper, we have presented some of these challenges. Also we have proposed a new research area of Cloud Control, which will address the whole range of closely related challenges for the efficient autonomic management of resources in cloud data centers.

References

- [1] J. Kephart and D. Chess. The vision of autonomic computing. *Computer*, 36(1):41–50, 2003.
- [2] M. Armbrust *et al.* A view of cloud computing. *Communications of the ACM*, 53(4):50–58, 2010.
- [3] J. Cao, M. Andersson, C. Nyberg, and M. Kihl. Web server performance modeling using an M/G/1/K*PS queue. In *10th International Conference on Telecommunications*, 2003.
- [4] D. Tipper and M. K. Sundareshan. Numerical methods for modeling computer networks under nonstationary conditions. *IEEE Journal on Selected Areas in Communications*, 8(9), 1990.
- [5] H. Khazaei, J. Mistic, and V. Mistic. Performance analysis of cloud computing centers using M/G/m/m+r queueing systems. *IEEE Transactions on Parallel and Distributed Systems*, Vol. 23, No. 5, May 2012.
- [6] K. Konstanteli, T. Varvarigou, and T. Cucinotta. Probabilistic admission control for elastic cloud computing. In *IEEE Int. Conf. on Service-Oriented Computing and Applications (SOCA)*, 2011, 2011.
- [7] W. Iqbala, M. Daileya, D. Carrerab, and P. Janecka. Adaptive resource provisioning for read intensive multi-tier applications in the cloud. *Future Generation Computer Systems*, Vol. 27, No. 6, 2010.
- [8] J. Tordsson, R. S. Montero, R.M. Vozmediano, and I.M. Llorente, Cloud brokering mechanisms for optimized placement of virtual machines across multiple providers, *Future Generation Computer Systems*, 28(2), 358-367, 2012.
- [9] A. Greenberg, J. Hamilton, D. Maltz, and P. Patel. The cost of a cloud: Research problems in data center networks. *ACM SIGCOMM Computer Communication Review*, 39(1), 2009.
- [10] M. Poess and R. Nambiar. Energy cost, the key challenge of today's data centers: A power consumption analysis of TPC-C results. *Proceedings of the VLDB Endowment*, 1(2), 2008.
- [11] E. Elmroth *et al.* . Self-Management Challenges for Multi-Cloud Architectures, *Towards a Service-Based Internet*, Lecture Notes in Computer Science, Vol. 6994, Springer-Verlag, 2011.
- [12] J. Heo, P. Jayachandran, I. Shin, D. Wang, T. Abdelzaher, and X. Liu. OptiTuner: On performance composition and server farm energy minimization application. *IEEE Transactions on Parallel and Distributed Systems*, Vol. 22, No. 11, 2011.