

# Automated Heterogeneous Service Management in Cloud Systems

Christopher C. Lamb, Jose-Marcio Luna, Pramod A. Jamkhedkar,

Gregory L. Heileman, Chaouki T. Abdallah

University of New Mexico

Department of Electrical and Computer Engineering

Albuquerque, NM 87131-0001

`{cclamb, jmarcio, pramod54, heileman, chaouki}@ece.unm.edu`

May 28, 2011

## Abstract

In this paper, we examine the problem of a single provider offering multiple types of service level agreements, and the implications thereof. In doing so, we propose a simple model for machine-readable service level agreements (SLAs) and outline specifically how these machine-readable SLAs can be constructed and injected into cloud infrastructures - important for next-generation cloud systems as well as customers. We then computationally characterize the problem, establishing the importance of both verification and solution, showing that in the general case injecting policies into cloud infrastructure is NP-Complete, though the problem can be made more tractable by further constraining SLA representations and using approximation techniques. Finally, we outline how to transform this problem into a more control-theoretic structure.

# 1 Introduction

The past few years have witnessed unprecedented expansion of commercial computing operations as the idea of cloud computing has become more mainstream and widely adopted by forward thinking technical organizational leadership. This rate of adoption promises to increase in the near future as well. With this expansion has come opportunity as well as risk, embodied by recent major service outages at leading cloud providers like Amazon. These issues promise to become more difficult to control as managed infrastructure expands. This expansion will simply not be possible without large amounts of automation in all aspects of cloud computing systems.

The current state of the art in cloud systems is poorly differentiated and not as customer-focused as it could be. Current providers place the responsibility of monitoring performance and proving outages on the consumer rather than providing more transparent and monitorable infrastructure [?]. Furthermore, providers as a whole usually provide one type of service level agreement (SLA) in a loosely-defined one-size-fits-all type of arrangement. This provides strong differentiating opportunities for smaller, second generation cloud system providers who have established the technology required to scalably manage multiple, competing SLAs on the same infrastructure in tandem with clear customer system visibility.

These second generation providers will rely on automated infrastructure management in order to scale. One of the first steps toward automating these systems is automating SLA management and compliance.

Herein, we will elaborate the idea of applying usage management to single system governed by multiple different types of SLAs. We will define the problem, more formally describe SLAs, analyze the implications of that formality, and using this information create a prototypical control system.

In Section 2, this paper begins by describing the different types of cloud computing models that generally exist today and how they manage services. Immediately thereafter, we propose a possible future model in which users can have unique SLAs that more closely fit their needs rather than shoehorning their computing needs into a previously configured contract. Then, in Section 3, we more formally define an SLA, and show how to convert one to an evaluable expression. In the following section, Section 4, we analyze the new SLA model and extract specific theoretical limits on computability and discuss implications thereof. Finally in Section ?? we use our new conclusions to design a prototypical control system using these principles.

## 1.1 Previous Work

As cloud computing is emerging as the future of utility systems hosting for consumer-facing applications. In these kinds of systems, components, applications, and hardware are provided as utilities over the Internet with associated pricing schemes pegged by system demand. Users accept specific QoS guidelines that providers use to provision and eventually allocate resources. These guidelines become the basis over which providers charge for services.

Over the past few years multiple service-based paradigms like web services, cluster computing and grid computing have contributed to the development of what we now call cloud computing [1]. Cloud computing distinctly differentiates itself from other service-based computing paradigms via a collective set of distinguishing characteristics: market orientation, virtualization, dynamic provisioning of resources, and service composition via multiple service providers [?]. This implies that in cloud computing, a cloud-service consumer's data and applications reside inside that cloud provider's infrastructure for a finite amount of time. Partitions of this data can in fact be handled by multiple cloud services,

and these partitions may be stored, processed and routed through geographically distributed cloud infrastructures. These activities occur within a cloud, giving the cloud consumer an impression of a single virtual system. These operational characteristics of cloud computing can raise concerns regarding the manner in which cloud consumer's data and applications are managed within a given cloud. Unlike other computing paradigms with a specific computing task focus, cloud systems enable cloud consumers to host entire applications on the cloud (i.e. Software as a Service) or to compose services from different providers to build a single system. As consumers aggressively start exploiting these advantages to transition IT services to external utility computing systems, the manner in which data and applications are handled within those systems by various cloud services will become a matter of serious concern.

A growing body of research has begun to appear over the past two years applying control theory to tuning computer systems. These range from controlling network infrastructure [2] to controlling virtualized infrastructure and specific computer systems [3], [4] to exploring feedforward solutions based on predictive modeling [5]. Significant open questions remain to research within this field [6], [7].

## 2 Cloud System Models

Current cloud systems do not ignore SLA restrictions; rather, they are designed from the ground up to support a single type of SLA. That SLA generally encompasses total system uptime and some kind of response time metric [?, ?]. If for some reason the cloud provider can no longer adhere to the terms outlined, some kind of compensation strategy applies to affected customers. Future cloud providers can very well use the ability to support multiple SLAs as a way to differentiate available products from competitors.



Figure 1: Amazon Auto Scaling


## 2.1 Current Model

Current systems like Amazon's EC2 or Rackspace products are designed around high availability, and this is reflected in the focus of their supplied SLAs. This common design focus is also evident in the artifacts generated by other vendors [?]. Furthermore, Amazon offers clear guidance on how to develop systems that take advantage of their robust architecture as well as services that provide some measure of automatic scaling [?, ?]. This combination of market leading position and products and the extensive supplied guidance make Amazon a clear choice to examine when reflecting on the current state-of-the-art.

Amazon's Cloud Watch products used in tandem with Auto Scaling provide the ability to control the number of deployed instances in response to specific system loads [?, ?]. Cloud Watch gives customers the ability to monitor various system performance metrics for their virtual machines, including but not limited to latency, processor use, and request counts. Furthermore, users can set resource levels at which additional EC2 instances are created or destroyed. This provides some level of personalized management and control over deployed systems within Amazon's cloud infrastructure.

## 2.2 Future Reference Model

While current cloud service providers focus on a single quality-of-service metric, future providers may very well begin to provide multiple metrics over which they will define service levels. This is not without precedent - just as airlines



Insert Image Here

Figure 2: Multiple SLA Architectural Integration

provide the same essential product at different service points, cloud providers could supply system hosting via disparate service levels, including divergent service metric definitions. For example, current architectures support uptime and availability as the primary managed metric from an SLA perspective. Future architectures could support uptime and availability, as well as specific latency, bandwidth, and geo-location sensitive hosting parameters. These kinds of SLAs would also continue to outline penalties when any of the conditions of that SLA were violated. Unlike current SLAs however, these could also differentiate based on the magnitude of the imposed penalty, with different classifications of service mapping to increasingly large penalties on service failure.

While industry does seem to certainly be trending in this direction, as indicated by the development of tools supporting user-centric infrastructure monitoring and management, this kind of control is not yet embedded into contracts of any kind, much less agreements that are machine-readable. Furthermore, this kind of management is still manual and cannot scale to the levels needed to manage Internet-scale systems.

### 3 Service Level Agreements Defined

As we have seen, SLAs generally consist of a set of conditions of use under which the SLA is binding, a set of obligations that the provider will adhere to if the customer adheres to the set conditions, and two sets of penalties, one penalizing

the provider when breaching obligations, and another penalizing the customer when breaching conditions of use. Conditions are generally loosely defined, while penalties are much more rigorously constructed. Generally however, conditions and obligations in this context can be viewed as defined by *objectives* which are measured by *indicators*. In the case of provider-centric obligations, these are commonly defined as Service Level Indicators (SLIs) and Objectives (SLOs).

With this general understanding of SLAs and related SLIs and SLOs, we can create a non-specific definition of an SLA as a set of tuples:

$$SLA = \{(I, O, E, P)_{0..n}\}, n \in Z \quad (1)$$

Where  $I$  is an indicator function,  $\forall i \in I, i : () \rightarrow \tau$ , which retrieve indicator values.  $O$  is a set of values,  $\forall o \in O, o : P(\tau)$ ,  $E$  is a set of predicates,  $\forall e \in E, e : ((() \rightarrow \tau) \times P(\tau) \rightarrow bool$ , and  $P$  is a set of penalty functions,  $\forall p \in P, p : T_{elapsed} \rightarrow Z$ .

For example, say we are a customer of Nimbus Cloud Corporation, and we have an SLA in which Nimbus provides guaranteed 100% uptime and packet latency between 300 and 750 milliseconds. If the first obligation is breached, Nimbus will pay me a fraction of my monthly fee depending on the length of the outage. If the second is breached, Nimbus will refund me \$0.01 per minute

of breach. This gives us a machine evaluateable SLA:

$$\begin{aligned}
SLA_{nimbus} = & ((uptime\_monitor() : bool, & (2) \\
& \{true\}, \\
& uptime\_evaluator(monitor : () \rightarrow bool, \{true\}) : bool, \\
& uptime\_penalty\_evaluator(T : Z) : Z), \\
& (latency\_monitor() : Z, \\
& \{300, 750\}, \\
& latency\_evaluator(monitor : () \rightarrow Z, \{300, 750\}) : bool, \\
& latency\_penalty\_evaluator(T : Z) : Z))
\end{aligned}$$

This more rigorous SLA allows users to monitor obligations and determine penalties when triggered.

## 4 Controlling with Service Level Agreements

Now that we have more rigorously defined our SLAs, notice that the SLA evaluation functions are predicates, and can be curried for later execution if needed. This allows us to begin a more fundamental analysis of SLAs and their capabilities.

### 4.1 Computational and Space Complexity

In Section 3, we defined an SLA to essentially be a sequence of evaluatable predicates. These evaluatable predicates are related in some way; currently, an SLA is the conjunction of these predicates. As these predicates can be created prior to evaluation and at evaluation time require no specific arguments once appropriately curried, we can define these predicates as as boolean *terms*. Ergo,



once we've created a group of predicates and transformed them into terms, we are evaluating an arbitrary boolean equation - in other words, we are verifying an instance of the *BooleanSatisfiabilityProblem*, or *SAT*.

*SAT* is *NP-Complete*, and very difficult to solve on today's computing systems [?]. *3-Sat*, a subset of *SAT*, is equally difficult, while *2-Sat* is not. *2-Sat* is firmly in the computational class *P*; in fact, *2-Sat* is *NL-Complete* as well, so we know it is solvable in an amount of space logarithmic in the number of boolean terms[?]. Furthermore, it is widely held that both *SAT* and *3-Sat* cannot be solved in logarithmic or less space.

## 4.2 Verification v. Solution

## 4.3 Approximation and other techniques

# 5 Control Theoretic Modeling

## 5.1 Results and Analysis

# 6 Conclusions and Future Works

## 6.1 Conclusions

## 6.2 Future Works

[8]

# References

- [1] R. Buyya, "Market-oriented cloud computing: Vision, hype, and reality of delivering computing as the 5th utility," in *Proceedings of the 2009 9th IEEE/ACM International Symposium on Cluster Computing and the Grid*,

- ser. CCGRID '09. Washington, DC, USA: IEEE Computer Society, 2009, pp. 1–. [Online]. Available: <http://dx.doi.org/10.1109/CCGRID.2009.97>
- [2] Y. Ariba, F. Gouaisbaut, and Y. Labit, “Feedback control for router management and tcp/ip network stability,” *Network and Service Management, IEEE Transactions on*, vol. 6, no. 4, pp. 255–266, 2009.
  - [3] Z. Wang, Y. Chen, D. Gmach, S. Singhal, B. Watson, W. Rivera, X. Zhu, and C. Hyser, “Appraise: application-level performance management in virtualized server environments,” *Network and Service Management, IEEE Transactions on*, vol. 6, no. 4, pp. 240–254, 2009.
  - [4] M. Kjaer, M. Kihl, and A. Robertsson, “Resource allocation and disturbance rejection in web servers using slas and virtualized servers,” *Network and Service Management, IEEE Transactions on*, vol. 6, no. 4, pp. 226–239, 2009.
  - [5] S. Abdelwahed, J. Bai, R. Su, and N. Kandasamy, “On the application of predictive control techniques for adaptive performance management of computing systems,” *Network and Service Management, IEEE Transactions on*, vol. 6, no. 4, pp. 212–225, 2009.
  - [6] X. Zhu, M. Uysal, Z. Wang, S. Singhal, A. Merchant, P. Padala, and K. Shin, “What does control theory bring to systems research?” *SIGOPS Oper. Syst. Rev.*, vol. 43, pp. 62–69, January 2009. [Online]. Available: <http://doi.acm.org/10.1145/1496909.1496922>
  - [7] J. Hellerstein, S. Singhal, and Q. Wang, “Research challenges in control engineering of computing systems,” *Network and Service Management, IEEE Transactions on*, vol. 6, no. 4, pp. 206–211, 2009.

- [8] Y. Chen, D. Gmach, C. Hyser, Z. Wang, C. Bash, C. Hoover, and S. Singhal, “Integrated management of application performance, power and cooling in data centers,” in *Network Operations and Management Symposium (NOMS), 2010 IEEE*, 2010, pp. 615 –622.