

Self-Optimizing Architecture for Ensuring Quality Attributes in the Cloud

Vivek Nallur, Rami Bahsoon
School of Computer Science, University of Birmingham
Birmingham B15 2TT, UK
{V.Nallur, R.Bahsoon}@cs.bham.ac.uk

Abstract

We describe various challenges in ensuring Quality Attributes(QA) of applications hosted in the cloud and hence the perceived quality of service of the cloud as a whole. We advocate a self-management/optimisation architecture-driven approach to ensure that Quality Attributes are met. We discuss the limitations of current approaches to self-managing architecture. We propose a novel approach, which exploits the El Farol problem as a modelling mechanism for QAs in architectures of applications in the cloud. The approach uses Service Level Agreements (SLA) and Utility Theory to direct the self-optimization. We conclude by looking at directions for further work.

1. Introduction

Advances in networking, storage and processing technologies have given us software that is mind-boggling in size and complexity of structure. Enterprises routinely deploy applications that span continents (e.g. through Grids, WANs, Internet), are composed of computing elements that are heterogeneous, and connected in complex topologies. The combination of development, configuration and management complexity has resulted in software fragility and increased interest in autonomic software [1], [2]. In the wake of organizational change, economic downturn and a demand for tightening the belt on IT costs, there is a trend toward moving large applications to the cloud [3]. Organizations such as IBM [4] and Gartner [5] advocate cloud computing as a potential cost-saver as well as provider of higher service quality.

Cloud computing refers to both the applications delivered as services over the Internet and the hardware and systems software in the datacenters that provide those services. [6]

Clouds, as made available by the major players like Amazon, Google and 3Tera, use the Software-As-A-Service or Infrastructure-As-A-Service model. This means that payment for the services of the cloud are made on the basis of cpu-hours used as well as storage used [7], which is highly economical. However cloud providers make no guarantees about the Quality of Service attributes being provided by

them. We envisage that Quality Attributes like performance, security, availability, reliability will be important parameters that organizations would want to monitor and optimize. It is in this context, that self-optimization of applications distributed in the cloud, becomes important.

Usually, applications in the cloud are of the following types:

- 1) Web-applications that cater to diverse users across the internet and demand fast response times.
- 2) Enterprise applications that cater to different business units across the world and require large amounts of secure, reliable data transfer and high availability (99.999%).
- 3) Scientific applications that need raw cpu or enterprises that perform batch processing

The third type of application has the natural advantage of 'cost associativity' offered by the cloud [6]. The first two types of applications represent a slice of systems that we depend on in our daily life; these need to be dependable and long-lived. By dependable, we refer to qualities like reliability, availability, and other non-functional requirements which need to be provided, maintained, evolved and monitored at runtime. These requirements often evolve in response to changes in users' needs, services, and the runtime environment. Consequently, this will necessitate dynamic change, evolution, and self-optimisation of applications hosted in the cloud with respect to the said requirements. Needless to say, evolution and self-optimisation needs to be constrained by cost, available resources, and any other operational constraints specific to the cloud.

The applicability of current research on self-managing architectures to the cloud is doubtful. This is because these approaches assume that the system is a closed-loop and self-management is done via a centralized component. As we will discuss in the section on Related Work, the environment of the cloud is sufficiently different (no control over physical topology, for instance), for centralized control and closed-loop type feedback to be inapplicable. The contribution of this paper is as follows: we highlight existing work on self-managing architectures and discusses their inapplicability to the case of the cloud. We turn to a motivating example to discuss the need for self-optimization in the cloud. We argue that Service Level Agreements (SLA) can assist the problem

of self-optimization in the cloud. We introduce the idea of an application changing/optimizing its architecture by means of modifying its components and connectors, through the use of composable web-services, that automatically negotiate their cost versus feature offering via SLAs. We propose to use the emergent effect of simple self-optimizations at lower levels (feature level) to achieve a higher level optimization (application-wide) in the cloud. We model the simple self-optimizations as solutions to the *El Farol Bar problem*.

The rest of this paper is organised as follows: In section 2, we present a motivating example to elucidate some of the Quality Attributes that need optimization. In section 3 we look at related work on self-managing architecture and discuss how it does not solve the problem mentioned in section 2. In section 4, we propose the idea of self-optimization in detail. And in the final section, we discuss future work.

2. Motivating Example

We motivate the need for research in self-managing architecture of cloud-based applications by the following example:

Consider a new social networking site (myChaseBook), much like Facebook, but in addition to the Web it also incorporates new media like cellphones and online radio. myChaseBook offers the following innovative services:

- Track friends geographically via cellphone, overlaid on a city map.
- Tagging of streaming radio.
- Recommend radio snippets to friends who can then chose to listen via cellphone and buy them from online stores.
- Cellphone-based access to photos, music tracks, friends and their recommendations based on the location of the user.

Suppose that myChaseBook becomes wildly popular and the creators of myChaseBook decide to move their site onto the Cloud. Their principal requirements are:

- **Scalability:** The application's design allows for scalability and hence the provisioning of computing resources like cpu, bandwidth should scale with demand. Cloud providers charge for their service on the basis of cpu-hours used, storage and bandwidth usage. Since myChaseBook is a startup and does not have too much cash in hand, the application should also scale down in times of low demand, to ensure maximum Return On Investment.
- **Availability:** Users are turned off if the service is not available due to frequent over-loading. The popularity of social networking sites tend to spike when websites like Slashdot recommend it or a celebrity starts using it. At such times, if the site is unavailable, then users tend

to turn away in search of alternatives. myChaseBook should be resilient in times of increased traffic and always be available.

- **QoS:** myChaseBook provides a lot of its services by composing several other web-services, that are created and controlled by third-parties. If these web-services fail to provide proper performance, then myChaseBook's reputation plummets. Hence, myChaseBook would like QoS from these web-services to be guaranteed via Service Level Agreements (SLA)
- **Cost:** The cost of using other web-services could fluctuate depending on demand and supply. myChaseBook should be able to automatically switch to a lower cost web-service, as long as it meets other requirements.

Using web-services in a composable manner to provide functionality is a typical instantiation of Service-Oriented Architecture. However for myChaseBook to consider using web-services in the cloud, its own design must have the following features:

- Service Discovery — myChaseBook should be able to discover other services that it can leverage to provide its own features.
- Contract negotiation for SLA [8] — myChaseBook's services should be able to automatically negotiate contracts with regard to performance and availability.
- Web-service provenance for process, data and workflow [9] — Payment for services could be based on data-flow or number of web-service invocations. In such cases, auditing of number of invocations and amount of data transfer would be necessary to ensure that correct payment is made.

Using these features, myChaseBook should be able to dynamically evolve its architecture by modifying the web-services to which it connects. When organizational priorities change, and myChaseBook wants to optimize on dependability (say), it could search for a web-service that offers an SLA with high dependability and use that web-service instead of the high performance web-service that it is currently using. However, decisions like changing focus from one Quality Attribute to another (performance to dependability) will probably not be optimal for the entire myChaseBook application. Each feature might want to optimize for a different Quality Attribute. For example, features like tagging of streaming radio and tracking of friends via cellphone might want to optimize for performance, while features like purchase of music tracks via cellphone might want to optimize for security. Further, all of these optimizations would have cost as a constraint as well. In this context, it is easy to see how different parts of an application might want to optimize differently, each realizing different tradeoffs.

3. Related Work

There has been a lot of work on self-managed architecture, as reported by Kramer and Magee [10]. Current research pursues software architectures as the appropriate level of abstraction for evaluating, reasoning about, managing and facilitating the dynamic change and evolution of complex software systems.

There have been a plethora of attempts at creating self-managing architectures. These attempts can be classified as approaching the problem from the following perspectives:

- Map system to ADL, dynamically change ADL, check for constraints, transform ADL to executable code [11], [12], [13].
- Create a framework for specifying types of adaptation possible using constraints, tactics thus ensuring ‘good adaptation’ [14], [15].
- Using formalized methods like SAM [16], control theory [17] and middleware optimizations [18].

However, all of these approaches envision a closed-loop system. That is, all of them assume that

- 1) The entire state of the application and the resources available for adaptation are known/visible to the management component.
- 2) The adaptation ordered by the management component is carried out in full, never pre-empted, ignored or disobeyed.
- 3) The management component gets full feedback or is able to view the results of changes made, on the entire system.

We feel that these are limiting assumptions because of the following reasons:

- 1) In a system of systems, each software component will have a different objective to attain and hence need not cooperate. In a cloud, we expect a system of systems to be implemented through composable web-services. Also, since each web-service is potentially owned or managed by a different entity, its optimization objective could be different.
- 2) In a truly large scale system, communication between centralised management components and distributed lower-level components may not be feasible/robust. For example, a performance optimizer component that is located on a server in the US may not be in a position to communicate adaptation orders to a component located on a server in Australia on a timely basis.
- 3) In a large scale system, global application knowledge is infeasible. If there are thousands of object instances, each with its own performance/memory issues, it is infeasible to expect a central optimizer to take into account each of those instances’ states or metrics while trying to optimize in reasonable time. This could be

due to multiple factors like combinatorial explosion, computational complexity or a fast rate of change.

4. Our Approach

4.1. Problem in Detail

It is easy to see that the architectural approaches proposed in the previous section do not fit the requirements of myChaseBook. The limitations elucidated preclude the achievement of the properties discussed in the motivating example. We envisage that applications that want to provide / consume services register themselves with the cloud, which acts as the central registry or even as a broker for negotiating price and SLAs. However, different web-services would have differing agendas and different properties that they wish to optimize for. It is in this context that this paper’s main contribution lies. We consider web-services running in the cloud ecosystem and consider how these services might optimize themselves and what this means for the quality of the cloud as a whole.

Each web-service will set its own priorities with regard to cost, performance, reliability and availability. Negotiation amongst these services will occur through any of a variety of resource reservation protocols as long as they conform to the *Web services agreement specification* [19]. The problem that we consider is whether such a system would reach its optimal level of cost vis-a-vis other quality attributes.

4.2. The El Farol solution

We view the distributed selfish self-optimization as the El Farol Bar problem [20] and seek inspiration in the solutions proposed in [21], [22]. Specifically, the constraint that we impose on any solution is that each web-service should attempt local self-management only, without any knowledge of the application as a whole. Also, it is not binding on any web-service to cooperate with any other web-service. *Yagan and Tham* describe a successful attempt [23] in using Reinforcement Learning to optimize network routing. We hope that a similar mechanism would work amongst web-services as well. A software architecture can be specified as a triad of components, connectors and topology [24]. In the particular case of the cloud, topology is opaque to the web-services and hence, the only optimization possible is in the types of components and connectors used.

4.2.1. Modelling. Modelling the web-services in a cloud in such a manner requires viewing the cloud as a marketplace. Each web-service can be viewed as an agent with a budget to spend and a targetted quality attribute to optimize. The cloud provider could be viewed as the ‘market maker’ or the broker, that specifies the negotiation protocol to be followed when web-services bid for services. Now, each agent applies

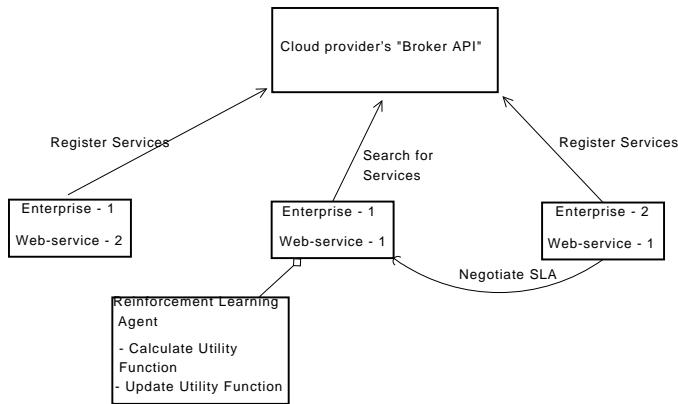


Figure 1. Interaction between autonomous web-services in the cloud

its own learning algorithm to choose amongst the strategies it has, to negotiate with other agents. Depending on the quality attribute that a particular agent is interested in optimizing, it would be prepared to pay a high price for an SLA that offers the desired guarantees. For example, an agent that is interested in optimizing performance would pay a premium for an SLA that guarantees a certain bandwidth or certain algorithmic performance. Alternatively, another agent interested in availability would pay a premium for an SLA that promises redundancy. How much premium to pay for obtaining a particular SLA is an engineering problem in decision making, with tradeoffs between multiple quality attributes and cost. A possible solution could be modelled using multi-attribute utility theory [25]. Allowing each agent to set its own parameters in a utility function creates a marketplace where each agent would set a (possibly) different price on different quality attributes. As shown in [22], if there are enough agents with differing strategies for maximising their utility, the system as a whole quickly reaches a near optimal level.

To return to the motivating example, myChaseBook's web-services would pay premiums for other web-services that provided (via SLAs) guarantees of performance and redundancy. We contend that with enough diversity in the population of web-services (and corresponding learning strategies), the cloud as a whole will reach a level of optimal cost vs. quality attribute tradeoff.

4.3. Applicability

Applications built in this manner should be resilient to changes in cloud conditions with regard to demand, change in types of supporting services and even organisational objectives. Thus, maintainability of applications in the cloud is greatly increased. For instance, a web-service that seeks to optimize on cost, would constantly monitor the demand being made on it and switch to cheaper composing services

as long as SLAs that it agreed to, are met. Now, if the organisation's objective were to change to providing a high level of availability, the web-service would (due to a change in its utility function) switch to more expensive services as long as they promise higher availability in their SLA.

5. Future Work

We think that modelling of the problem in greater detail would be necessary to carry out any experimental evaluation. We will propose more precise application of multi-attribute utility theory to SLA negotiation. Also, modification strategies that a web-service can learn from, will be proposed. This would enable simulations of a cloud with negotiating web-services, thus allowing us to test our idea of low-level self-optimization leading to an emergent higher level optimized application state in the cloud. If successful, this would lead to long-lived applications in the cloud being more resilient to change, and successfully adapting to changing Quality Attribute optimization needs.

References

- [1] M. Parashar and S. Hariri, *Autonomic Computing: An Overview*, 2005. [Online]. Available: http://dx.doi.org/10.1007/11527800_20
- [2] IBM, "An architectural blueprint for autonomic computing," June 2006. [Online]. Available: http://www-01.ibm.com/software/tivoli/autonomic/pdfs/AC_Blueprint_White
- [3] C. Boulton, *Forrester's Advice to CFOs: Embrace Cloud Computing to Cut Costs*. eWeek.com, October 2008. [Online]. Available: <http://www.eweek.com/c/a/Enterprise-Applications/Forresters-Advice-to-CFO>
- [4] IBM, "Ibm perspective on cloud computing," Whitepaper published on the web, <http://www-05.ibm.com/services/ch/download/ibm-cloud-computing-12-08.pdf>, 2008. [Online]. Available: <http://www.ibm.com/ibm/cloud/>
- [5] C. Boulton, "Gartner sees great saas enterprise app growth despite downturn," eWeek.com, October 2008. [Online]. Available: <http://www.eweek.com/c/a/Enterprise-Applications/Gartner-Sees-Great-SAA>
- [6] M. Armbrust, A. Fox, R. Griffith, A. D. Joseph, R. H. Katz, A. Konwinski, G. Lee, D. A. Patterson, A. Rabkin, I. Stoica, and M. Zaharia, "Above the clouds: A Berkeley view of cloud computing," EECS Department, University of California, Berkeley, Tech. Rep. UCB/EECS-2009-28, Feb 2009. [Online]. Available: <http://www.eecs.berkeley.edu/Pubs/TechRpts/2009/EECS-2009-28.html>
- [7] A. Inc, *Amazon Elastic Compute Cloud (Amazon EC2)*. <http://aws.amazon.com/ec2/#pricing>: Amazon Inc., 2008. [Online]. Available: <http://aws.amazon.com/ec2/#pricing>

- [8] R. Buyya, C. S. Yeo, and S. Venugopal, "Market-Oriented cloud computing: Vision, hype, and reality for delivering IT services as computing utilities," *0808.3558*, August 2008. [Online]. Available: <http://arxiv.org/abs/0808.3558>
- [9] B. Hayes, "Cloud computing," *Commun. ACM*, vol. 51, no. 7, pp. 9–11, 2008. [Online]. Available: <http://dx.doi.org/10.1145/1364782.1364786>
- [10] J. Kramer and J. Magee, "Self-managed systems: an architectural challenge," in *FOSE '07: 2007 Future of Software Engineering*. Washington, DC, USA: IEEE Computer Society, 2007, pp. 259–268. [Online]. Available: <http://dx.doi.org/10.1109/FOSE.2007.19>
- [11] E. M. Dashofy, A. van der Hoek, and R. N. Taylor, "Towards architecture-based self-healing systems," in *WOSS '02: Proceedings of the first workshop on Self-healing systems*. New York, NY, USA: ACM Press, 2002, pp. 21–26. [Online]. Available: <http://dx.doi.org/10.1145/582128.582133>
- [12] —, "An infrastructure for the rapid development of XML-based architecture description languages," in *ICSE '02: Proceedings of the 24th International Conference on Software Engineering*. New York, NY, USA: ACM Press, 2002. [Online]. Available: <http://dx.doi.org/http://dx.doi.org/10.1145/581339.581374>
- [13] —, "A comprehensive approach for the development of modular software architecture description languages," *ACM Trans. Softw. Eng. Methodol.*, vol. 14, no. 2, pp. 199–245, April 2005. [Online]. Available: <http://dx.doi.org/10.1145/1061254.1061258>
- [14] I. Georgiadis, J. Magee, and J. Kramer, "Self-organising software architectures for distributed systems," in *WOSS '02: Proceedings of the first workshop on Self-healing systems*. New York, NY, USA: ACM Press, 2002, pp. 33–38. [Online]. Available: <http://dx.doi.org/10.1145/582128.582135>
- [15] s. Cheng, D. Garlan, and B. Schmerl, "Architecture-based self-adaptation in the presence of multiple objectives," in *Proceedings of the 2006 international workshop on Self-adaptation and self-managing systems*. Shanghai, China: ACM, 2006, pp. 2–8. [Online]. Available: <http://dx.doi.org/10.1145/1137677.1137679>
- [16] J. Wang, C. Guo, and F. Liu, "Self-healing based software architecture modeling and analysis through a case study," in *Networking, Sensing and Control, 2005. Proceedings. 2005 IEEE*, 2005, pp. 873–877. [Online]. Available: <http://dx.doi.org/10.1109/ICNSC.2005.1461307>
- [17] J. Hellerstein, *Engineering Self-Organizing Systems*, 2007, p. 1. [Online]. Available: http://dx.doi.org/10.1007/978-3-540-74917-2_1
- [18] M. Trofin and J. Murphy, "A Self-Optimizing container design for enterprise java beans applications," in *In Proceedings of the Second International Workshop on Dynamic Analysis (WODA 2004, 2003*. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.59.2979>
- [19] A. Andrieux, K. Czajkowski, A. Dan, K. Keahey, H. Ludwig, T. Nakata, J. Pruyne, J. Rofrano, S. Tuecke, and M. Xu, "Web services agreement specification," <http://xml.coverpages.org/WS-Agreement-13652.pdf>. [Online]. Available: <http://xml.coverpages.org/WS-Agreement-13652.pdf>
- [20] B. W. Arthur, "Inductive reasoning and bounded rationality," *American Economic Review (Papers and Proceedings)*, no. 84, pp. 406–411, 1994. [Online]. Available: http://www.santafe.edu/~wbarthur/Papers/EL_Farol.html
- [21] D. H. Wolpert and K. Tumer, *An Introduction to Collective Intelligence*, August 1999. [Online]. Available: <http://arxiv.org/abs/cs.LG/9908014>
- [22] D. H. Wolpert, K. R. Wheeler, and K. Tumer, "Collective intelligence for control of distributed dynamical systems," *Europhysics Letters*, vol. 49, pp. 708–714, 2000. [Online]. Available: <http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.43.3542>
- [23] D. Yagan and C. K. Tham, "Self-optimizing architecture for qos provisioning in differentiated services," in *ICAC '05: Proceedings of the Second International Conference on Automatic Computing*. Washington, DC, USA: IEEE Computer Society, 2005, pp. 346–347. [Online]. Available: <http://dx.doi.org/10.1109/ICAC.2005.57>
- [24] L. Bass, P. Clements, and R. Kazman, *Software Architecture in Practice, Second Edition*, 2nd ed. Addison-Wesley Professional, April 2003. [Online]. Available: <http://my.safaribooksonline.com/0201703726>
- [25] D. L. Thurston, *Decision Making in Engineering Design*. ASME Press. New York, New York, 2006, ch. Multi-attribute Utility Analysis of Conflicting Preferences, pp. 125–133.