

POLITECNICO DI MILANO

Scuola di Ingegneria dell'Informazione
Corso di Laurea Magistrale in Ingegneria Informatica
Dipartimento di Elettronica, Informazione e Bioingegneria



Generalized Nash Equilibria for the Service Provisioning Problem in Multi-Cloud Systems

Advisor: Prof. Danilo ARDAGNA

Co-Advisor: Prof. Mauro PASSACANTANDO

Master Thesis by:
Ettore TREVISIOL - 770628

Academic Year 2012/2013

*This thesis is dedicated to my parents.
For their endless love, support and encouragement.*

Acknowledgements

I would like to express my deepest gratitude to my advisor, Prof. Danilo Ardagna, for his excellent guidance, caring, patience and providing me with an excellent atmosphere for doing this thesis. My thanks go also to Prof. Mauro Passacantando for his great help in mathematical issues while developing this work.

I wish to thank my university mates, in particular thanks go to Giorgio Spadaro, Paolo Salvini and Riccardo Sacchi. I am also grateful to Davide Molinari, Anna Savi and the other students of the Software Engineering Laboratory for sharing their knowledge with me during the development of this thesis.

Finally, a particular thank goes to Andrea Fresco and Pier Paolo Cedaro, friends with whom I shared great moments of life.

Ettore

Contents

1	Introduction	5
2	State of the art	9
2.1	Cloud Computing basic concepts	9
2.2	Cloud Computing definition	11
2.2.1	Characteristics	13
2.2.2	Structure models	14
2.3	Cloud Computing and run-time research challenges	23
2.3.1	Problem	25
2.3.2	Solution	25
2.3.3	Discipline	26
2.3.4	State of the art	27
2.3.5	Classification of the state of the art	33
2.3.6	Criteria for evaluation	38
3	A game theory service provisioning model	41
3.1	Problem statement and assumptions	41
3.2	Generalized Nash game model	44
3.3	Game analysis	48
3.3.1	Dominant strategies for IaaS	48
3.3.2	Game potential	48
3.3.3	Analysis of constraints	49
3.3.4	Game model reformulation	52
3.4	A distributed algorithm for identifying Generalized Nash Equilibria	53
4	Tools	59
4.1	AMPL	59
4.2	CPLEX	60
4.2.1	CPLEX algorithms for continuous optimization	60
4.3	SPECweb 2005	61

CONTENTS

4.4	JMeter	63
4.5	SPECweb deployment in the Cloud	65
4.5.1	SPECweb tests	65
4.5.2	JMeter extension	67
4.5.3	SPECmeter	69
4.6	Cloud analysis tool	71
4.6.1	Cloud analysis tool class diagram	71
4.6.2	Cloud analysis tool sequence diagram	72
5	Experimental results	75
5.1	Design of experiments	75
5.1.1	Parameters generation	75
5.1.2	SaaS to IaaS mapping	77
5.1.3	Traffic generation	77
5.2	Scalability analysis	80
5.3	Equilibria efficiency	81
5.3.1	Social optimum problem	82
5.3.2	Price of Anarchy and Individual Worst Case	83
5.4	Alternative algorithms for resource allocation	83
5.4.1	Heuristic	84
5.4.2	Resource rescaling algorithm	89
5.5	Algorithms efficiency comparison	90
5.6	Multiple IaaS analysis	96
6	Conclusions	105
	Appendices	109
A	Game theory and generalized Nash equilibrium problem	109
A.1	Game theory in the Cloud Computing	109
A.2	Definition of Game	110
A.3	Solution concepts: Nash Equilibrium and Generalized Nash Equilibrium	111
A.4	Equilibria existence and potential games	114
A.5	Wardrop equilibrium	117

List of Figures

2.1	Cloud Computing architecture, [106].	15
2.2	Cloud service models.	18
2.3	Cloud deployment models.	22
2.4	Taxonomy for optimization approaches.	28
3.1	Cloud infrastructures.	42
3.2	System performance model.	43
3.3	Algorithm for finding Generalized Nash Equilibria.	55
4.1	SPECweb 2005 test diagram.	66
4.2	SPECweb banking test Markov chain.	67
4.3	Markov chain example in JMeter.	68
4.4	SPECmeter architecture.	69
4.5	SPECmeter test automation sequence diagram.	70
4.6	Cloud analysis tool class diagram.	72
4.7	Cloud analysis tool sequence diagram.	73
5.1	Queueing delay time.	76
5.2	Service time.	76
5.3	Daily time distribution of requests.	78
5.4	Weekly time distribution of requests.	79
5.5	Generated daily time distribution of requests.	79
5.6	Worldwide distribution of requests.	80
5.7	Distributed algorithm for identifying GNE scalability.	82
5.8	Impact of system capacity and service classes on the optimal profit, [59].	84
5.9	The difference of optimal profit with varied number of opened service classes, [59].	85
5.10	IaaS capacity usage on peak hours.	92
5.11	Traffic example of a single SaaS provider.	97
5.12	SaaS allocation with unlimited resources.	98

LIST OF FIGURES

5.13 SaaS allocation with limited resources.	99
5.14 Multi-IaaS analysis results with $\phi_i = 0.1$	100
5.15 Multi-IaaS objective function comparison with $\phi_i = 0.1$	101
5.16 Multi-IaaS analysis results with $\phi_i = 0.3$	102
5.17 Multi-IaaS objective function comparison with $\phi_i = 0.3$	103
5.18 Multi-IaaS analysis results with $\phi_i = 0.5$	103
5.19 Multi-IaaS objective function comparison with $\phi_i = 0.5$	104
A.1 Families of Generalized Nash Equilibrium Problems.	114

List of Tables

2.1	Problem category: perspective.	33
2.2	Problem category: quality attributes.	34
2.3	Problem category: dimensionality.	34
2.4	Problem category: constraints.	35
2.5	Solution category: type.	35
2.6	Solution category: degrees of freedom.	36
2.7	Solution category: architecture representation.	36
2.8	Solution category: optimization strategy.	36
2.9	Solution category: constraints handling.	37
2.10	Solution category: time scale.	37
2.11	Discipline category: type.	37
2.12	Discipline category: quality model.	37
3.1	Parameters and decision variables.	47
5.1	Performance parameters and time unit costs.	77
5.2	Distributed algorithm for identifying GNE execution times. . .	81
5.3	Algorithms PoA comparison with $\phi_i = 0.1$	93
5.4	Algorithms IWC comparison with $\phi_i = 0.1$	93
5.5	Algorithms re-optimization comparison with $\phi_i = 0.1$	93
5.6	Algorithms PoA comparison with $\phi_i = 0.3$	94
5.7	Algorithms IWC comparison with $\phi_i = 0.3$	94
5.8	Algorithms re-optimization comparison with $\phi_i = 0.3$	94
5.9	Algorithms PoA comparison with $\phi_i = 0.5$	95
5.10	Algorithms IWC comparison with $\phi_i = 0.5$	95
5.11	Algorithms re-optimization comparison with $\phi_i = 0.5$	95
5.12	Multi-IaaS analysis results with $\phi_i = 0.1$	99
5.13	Multi-IaaS analysis results with $\phi_i = 0.3$	100
5.14	Multi-IaaS analysis results with $\phi_i = 0.5$	101

LIST OF TABLES

Abstract

In recent years the evolution and the widespread adoption of virtualization, service-oriented architectures, autonomic, and utility computing have converged letting a new paradigm to emerge: Cloud Computing. Clouds allow the on-demand delivering of software, hardware, and data as services. Currently the Cloud offer is becoming day by day wider, since all the major IT companies and Service providers, like Microsoft, Google, Amazon, HP, IBM, and VMware have started providing solutions involving this new technological paradigm.

As Cloud-based services are more numerous and dynamic, the development of efficient service provisioning policies becomes increasingly challenging. In this thesis we take the perspective of Software as a Service (SaaS) providers which host their applications at multiple Infrastructure as a Service (IaaS) providers. Each SaaS needs to comply with quality of service requirements, specified in Service Level Agreement (SLA) contracts with the end-users, which determine the revenues and penalties on the basis of the achieved performance level. Each SaaS provider wants to minimize the cost of use of Cloud resources and penalties for requests execution failures. Moreover, SaaS providers compete and bid for the use of infrastructural resources. On the other hand, the IaaS want to maximize their revenues obtained providing virtualized resources.

In this thesis we model the service provisioning problem as a generalized Nash game. In particular, we develop an efficient distributed algorithm for the run-time allocation of IaaS resources among competing SaaS providers. We demonstrate the effectiveness of our approach by performing tests considering realistic Cloud scenarios. Numerical results show that our algorithm is scalable and can be used at run-time since it can solve problem instances of maximum size in less than one minute. Compared to other state-of-the-art solutions our model can improve the efficiency of Cloud system evaluated in term of Price of Anarchy up to 100%. Furthermore our analyses point out the SaaS benefits while exploiting multiple IaaS deployment of applications and redistribution of application workloads.

Sommario

Negli ultimi anni l'evoluzione e la diffusa adozione di virtualizzazione, architetture orientate ai servizi, autonomic e utility computing sono confluiti in un nuovo paradigma: il Cloud Computing. Il Cloud Computing ha come obiettivo la fornitura on-demand di software e hardware come risorse accessibili tramite Internet.

Con l'aumento della quantità e della dinamicità dei servizi basati sul Cloud, lo sviluppo di politiche efficienti per la distribuzione delle risorse è diventato sempre più complesso. In questo lavoro di tesi abbiamo studiato il problema dal punto di vista dei fornitori Software as a Service (SaaS) che ospitano le loro applicazioni presso molteplici fornitori Infrastructure as a Service (IaaS). Ogni SaaS deve rispettare la qualità del servizio, specificata nei contratti di Service Level Agreement (SLA) con i propri clienti, che determina i ricavi e le penalità sulla base del livello di prestazioni raggiunto. Ogni SaaS vuole minimizzare i costi di utilizzo delle risorse Cloud e delle sanzioni causate dalla violazione dei contratti di SLA. Inoltre, i fornitori SaaS competono fra loro facendo offerte per l'utilizzo delle infrastrutture. Al contrario, gli IaaS vogliono massimizzare i propri introiti ottenuti fornendo risorse virtualizzate.

In questa tesi abbiamo modellato il problema per la fornitura delle risorse come un gioco di Nash generalizzato. In particolare, abbiamo sviluppato un algoritmo distribuito per la gestione a run-time delle risorse degli IaaS fra i SaaS in competizione. Abbiamo dimostrato l'efficacia del nostro approccio compiendo test rappresentativi di scenari di carico reale. I risultati numerici mostrano che l'algoritmo è scalabile e può essere utilizzato a run-time, poiché può risolvere istanze del problema di dimensione massima in meno di un minuto. Rispetto ad altre soluzioni della letteratura il nostro modello può migliorare l'efficienza in termini di Price of Anarchy del sistema Cloud valutato fino al 100%. Inoltre, le nostre analisi evidenziano i benefici che i SaaS possono ottenere sfruttando il dislocamento delle applicazioni e la distribuzione del traffico su molteplici IaaS.

Chapter 1

Introduction

Cloud Computing has been a dominant IT news topic over the past few years. It is essentially a way for IT companies to deliver software/hardware on-demand as services through the Internet. Cloud Computing applications are generally priced on a subscription model, so end-users may pay a yearly usage fee, for example, rather than the more familiar model of purchasing software licenses. The Cloud-based services are not only restricted to software applications (Software as a Service – SaaS), but could also be the platform for the deployment and execution of applications developed in house (Platform as a Service – PaaS) and the hardware infrastructure (Infrastructure as a Service – IaaS).

In the SaaS paradigm, applications are available over the Web and provide Quality of Service (QoS) guarantees to end-users. The SaaS provider hosts both the application and the data, hence the end-user is able to use and access the service from all over the world. With PaaS, applications are developed and deployed on platforms transparently managed by the Cloud provider. The platform typically includes databases, middleware, and also development tools. In IaaS systems, virtual computer environments are provided as services and servers, storage, and network equipment can be outsourced by customers without the expertise to operate them.

Many companies, e.g., Google, Amazon, and Microsoft are offering Cloud Computing services such as Google’s App Engine and Amazon’s Elastic Compute Cloud (EC2) or Microsoft Windows Azure. Large data centers provide the infrastructure behind the Cloud and virtualization technology makes Cloud Computing resources more efficient and cost-effective both for providers and customers. Indeed, end-users obtain the benefits of the infrastructure without the need to implement and administer it directly adding or removing capacity almost instantaneously on a “pay-as-you-use” basis. Cloud providers can, on the other hand, maximize the utilization of their

physical resources also obtaining economies of scale.

The development of efficient service provisioning policies is among the major issues in Cloud research. Indeed, modern Clouds operate in a new and dynamic world, characterized by continuous changes in the environment and in the system and performance requirements must be satisfied. Continuous changes occurs without warning and in an unpredictable manner, and are outside the control of the Cloud provider. Therefore, advanced solution need to be developed to manage the Cloud system in a dynamically adaptive way, while continuously providing service and performance guarantees.

The recent evolution of Cloud system and the rapid growth of the Internet have led to a remarkable usage of game-theoretic tools. Problems arising in the ICT industry, such as resource or quality of service allocation problem, pricing, and load shedding, can not always be handled with classical optimization approaches because each player can be affected by the action of all players, not only by his own actions. In this context, game theory models and approaches allow to gain an in-depth analytical understanding of the service provisioning problem.

Game Theory has been successfully applied to diverse problems such as Internet pricing, flow and congestion control, routing, and networking. One of the most widely used “solution concept” in Game Theory is the Nash Equilibrium approach [71]: a set of strategies for the players constitute a Nash Equilibrium if no player can benefit by changing his/her strategy unilaterally or, in other words, every player is playing a *best response* to the strategy choices of his/her opponents.

In this thesis we take the perspective of SaaS providers which host their applications at multiple IaaS providers, thanks to a software layer developed within the MODAClouds project [68] [13]. Each SaaS provider wants to minimize the cost of use of Cloud resources and incurs in penalties in case of requests execution failures. The cost minimization is challenging since online services receive dynamic workloads that fluctuate during the day. Resources have to be allocated flexibly at run-time according to workload fluctuations. Furthermore, each SaaS behaves selfishly and competes with others SaaS for the use of infrastructural resources supplied by the IaaS. Each IaaS, in his turn, wants to maximize the revenues obtained providing the resources.

To capture the behavior of SaaSs and IaaS in this conflicting situation in which the best choice for one depends on the choices of the others, we recur to the *Generalized Nash Equilibrium* (GNE) concept, which is an extension of the classical Nash equilibrium.

Therefore, the run-time service provisioning problem will be modeled as a *Generalized Nash Equilibrium Problem* (GNEP). We then use Noncooperative Game Theory results to develop an efficient algorithm for the run-time

management and allocation of IaaS resources to competing SaaS suitable also for a *fully distributed* implementation. Multiple solutions achieving generalized equilibria will be proposed and evaluated in terms of their efficiency with respect to the *social optimum* of the Cloud. We will demonstrate the effectiveness of our approach by simulation and performing tests on a real prototype environment.

The remainder of the thesis is organized as follows:

- In Chapter 2 we present a general overview on Cloud Computing providing definitions, illustrating the main characteristics and showing the different structures models available. Afterwards, we will explain the state of the art concepts and techniques relative to our work. An analysis and a classification of the literature approaches is given in terms of type of problem, solution found and discipline adopted, according to the approach used: pure optimization or game theory.
- In Chapter 3 the game-theoretic service provisioning problem will be faced. We will start introducing the problem statements and design assumptions. Then an analysis of the game and its properties will be performed. Finally, we will present an algorithm able to find an equilibrium for the resource allocation problem.
- In Chapter 4 we will describe the tools used in this thesis work. A description of the optimization modelers and solvers we adopted will be provided. Furthermore, we will describe the workload injector we have developed to estimate the performance parameters of Cloud applications.
- Chapter 5 will be dedicated to assess the quality of our solution through analyses and experiments. After a description of experiments settings, algorithm scalability and comparison with other two approaches will be performed. Finally, we will analyze the benefits that can be achieved by SaaS when hosting applications on multiple IaaS.
- In Chapter 6 conclusions will be drawn, underling the achieved results and presenting future research directions.

Chapter 6

Conclusions

As Cloud-based services become more numerous and dynamic, resource provisioning becomes more and more challenging, especially when decisions can be affected by the action of others, not only by own actions, hence requiring game-theory approaches. Indeed, in any time instant resources have to be allocated to handle effectively workload fluctuations, while providing quality of service guarantees to the end users.

The overall goal we addressed in our thesis is the minimization of the costs associated with the allocated virtual machine instances in multiple IaaS, while guaranteeing QoS constraints. To achieve this purpose we have proposed a game-theoretic approach for the run-time management of IaaS provider capacities among multiple competing SaaS. The cost model consists of objective functions which include revenues and penalties incurred depending on the achieved performance level and infrastructural costs associated with IaaS resources. Therefore a distributed algorithm for identifying Generalized Nash Equilibria have been presented and its termination in a finite number of iterations has been demonstrated.

Thanks the AMPL language and the CPLEX solver, the effectiveness of our approach have been assessed by performing a wide set of analyses under multiple workload conditions. Realistic workloads created from a large website statistics and performance parameters estimated on an industrial benchmarking deployed in the Cloud have been used.

A number of different scenario of interest have been considered. Systems up to thousands of applications can be managed very efficiently in a fully distributed manner. Our algorithm found efficient GNE, for a hourly basis resource allocation, in less than a minute, proving to be perfectly suitable for run-time provisioning.

A comparison with utilization based state-of-the-art techniques and a rescaling algorithm shows that our solution outperform alternative methods

proving better results in terms of equilibrium efficiency both as regards the PoA both the IWC of each SaaS. In addition, our algorithm, achieve an efficient GNE under heavy workload conditions while thresholds based heuristic finds infeasible results due to the inability to manage SaaSs competition.

Finally analyses showed clear SaaS benefits while exploiting multiple IaaS deployment of applications and redistribution of traffic. SaaSs can have an average savings up to 50% compared to single IaaS architectures. The equilibria achieved are close to Cloud optimum, with inefficiency less than 2% for the social problem and of 6% for individual player in the worst cases.

Future work will extend the proposed solutions to consider multiple time-scales for performing resource allocation from few minutes to one hour. Additionally short-term solutions will be based on receding horizon techniques.