# Anticipatory Models of Load Balancing in Cloud Computing

A Thesis Submitted for the Degree of

Doctor of Philosophy

By

Shahrzad Aslanzadeh

**UNIVERSITY OF TECHNOLOGY SYDNEY**

University of Technology Sydney

New South Wales, Australia

CERTIFICATE OF ORIGINAL AUTHORSHIP

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Signature of Student: *Shahrzad Aslanzadeh*

Date: 04/14/2016

# Acknowledgment

I would like to express my sincere gratitude to my supervisor Dr. Zenon Chaczko for his support, continuous guidance, meticulous suggestions and his inexhaustible patience especially during preparation of this dissertation. At many stages in the course of this research, I benefited from his advice, particularly so when exploring new ideas.

Also I would like to thank Dr. Christopher Chiu, who as a good friend was always willing to help and give his best suggestions. Many thanks to Christopher Mcdermid for helping me through my case study design.

I will be thankful to my close friends for their friendly support and caring.

I express my personal love and appreciation to my mum, dad and my sister for constant motivation and support.

Finally, I would like to thank my husband Alireza, who was always there cheering me up and stood by me through the good times and the bad.

# Abstract

Cloud Computing is a recent arrival to the world of IT infrastructure. The concept allows companies to maximise utilisation of their potentials and consequently boost their performance. One of the main benefits of Cloud Computing is the significant increase in efficiency of executing business plans. Additionally, Cloud Computing provides large-scale applications with powerful computing power across global locations. Yet Cloud users are able to share their data easily by using replication methodologies.

Cloud Computing structure has been developed based on a multi-tenancy concept. Therefore, availability and efficiency of the resources are important factors in the Cloud architecture. However, as the numbers of users are increasing rapidly, the load will have a significant impact on performance and operation of the Cloud systems. Accordingly, optimised load balancing algorithms that can manage the Cloud load in a time- and cost-efficient manner are required.

Much research in recent years has been dedicated to optimising load balancing in Cloud Computing. This optimisation is demonstrated through a balanced network of interacting resources. The goal of this network is to minimise the wait time and maximise utilisation of the throughput.

This thesis provides a set of solutions which mitigate the problem of load balancing in the Cloud. The dissertation investigates a novel class of heuristic scheduling algorithms that improves load balancing in workflow scheduling applications.
Furthermore, it proposes a new anticipatory replication methodology with the objective of improving data availability to enhance the load balancing between the Cloud sites.

In summary, this research innovation implicates the design of optimised load balancing algorithms that consider the magnitude and direction of the load in workflow applications. Furthermore, by architecting the anticipatory replication algorithm, it minimises the numbers of the replicas and enhances the effective network usage in Cloud-based systems.

Contents

# List of Figures

# List of Tables

# Glossary

| | |
|---|---|
| AMAZON EC2 | Amazon Elastic Cloud Computing |
| ANM | Anisotropic Network Model |
| DAG | Directed Acyclic Graph |
| DPSO | Discrete Version Of PSO |
| ENM | Elastic Network Models |
| FCFS | First Come First Served |
| FMOC | Finite Multi-Order Context |
| GA | Genetic Algorithm |
| GBEST | Global Best Position |
| GNM | Gaussian Network Model |
| GRMS | Global Replica Management System |
| HDM | Hospital Data Management |
| HEFT | Heterogeneous Earliest Finish Time |
| HTML | Hyper Text Mark-up Language |
| I/O | Input/Output |
| IAAS | Infrastructure As A Service |
| ICT | Information & Communication Technology |
| IOT | Internet Of Things |
| IT | Information Technology |
| LFU | Least Frequently Used |
| LRU | Least Recently Used |
| MAKESPAN | Total Length Of Schedule |
| MBIT/S | Megabit Per Second |
| NIST | National Institute Of Standards And Technology |
| NP PROBLEM | Non-Deterministic Polynomial Time |
| PAAS | Platform As A Service |
| PBEST | Best Local Position |
| PSO | Particle Swarm Optimisation |
| QOS | Quality Of Service |
| RC | Relative Cost |

| | |
|---|---|
| SAAS | Software As A Service |
| SDDRC | Smart Dynamic Data Replication In Cloud Computing |
| SHEFT | Scalable Heterogeneous Earliest Finish Time |
| SLA | Service Level Agreement |
| STEM | Generalised Spring Tensor Model |
| STEM-PSO | Generalised Spring Tensor Model- Particle Swarm Optimisation |
| URL | Uniform Resource Locator |
| VM | Virtual Machine |
| XML | Extensible Mark-up Language |
| XMPP | Extensible Messaging And Presence Protocol |