

## **UNIVERSITI PUTRA MALAYSIA**

SCHEDULING TIGHT DEADLINES FOR SCIENTIFIC WORKFLOWS IN THE CLOUD

AWADH SALEM SALEH BAJAHER

FSKTM 2018 54



## SCHEDULING TIGHT DEADLINES FOR SCIENTIFIC WORKFLOWS IN THE

CLOUD

By

## AWADH SALEM SALEH BAJAHER

Thesis Submitted to the School of Graduate Studies, University Putra Malaysia, in Fulfillment of the Requirement for the Degree of Master of Computer science

JULY 2018

## COPYRIGHT

All material contained within the thesis, including without limitation text, logos, icons, photographs and all other artwork, is copyright material of University Putra Malaysia unless otherwise stated. Use may be made of any material contained within the thesis for non-commercial purposes from the copyright holder. Commercial use of material may only be made with the express, prior, written permission of University Putra Malaysia.

Copyright ©University Putra Malaysia

Abstract of thesis presented to the Senate of University Putra Malaysia in Fulfilment of the Requirement for the Degree of Master of Computer

Science

## SCHEDULING TIGHT DEADLINES FOR SCIENTIFIC WORKFLOWS IN

#### THE CLOUD

By

AWADH SALEM SALEH BAJAHER

**JULY 2018** 

Chair: NOR ASILA WATI ABDUL HAMID, PhD Faculty: Computer Science and Information Technology

#### Abstract

Cloud computing has increasingly become a demand for scientific computations as it provides users with simple access for computation. Commercial clouds are also used for scientific analysis and computation because of their scalability, latest high-quality hardware as well as pay-peruse cost model. Commercial clouds can be easily accessed globally. There have been several studies presenting new algorithms to generate deadline constrained schedules to minimize the execution cost as well as the high failure rate in schedule constructions. However, there are increased failure rates whenever tight deadlines are produced.

The work in this paper focuses on the hurdle of scheduling tight deadline scientific workload. This article will evaluate the performance of the

Proportional Deadline Constrained (PDC) algorithm using Cloudsim and compare it with the Deadline Constrained Critical Path (DCCP) scheduling algorithm. The performance evaluation is done using two different performance metrics, success rate and normalized cost. The results show that the PDC performs better in term of success rate metric while the DCCP algorithm has better performance in term of normalized cost metric. The PDC could be improved on the normalized cost.

**Keywords** Deadline constrained, scientific workflows scheduling, cloud resources, resource provisions.

## **APPROVAL FORM**

This thesis was submitted to the Senate of University Putra Malaysia and has been accepted as fulfilment of the requirement for the degree of Master of Computer science. The members of the Supervisory Committee were as follows:

Nor Asila Wati Abdul Hamid, PhD Faculty of computer science and IT technology University Putra Malaysia (Supervisor)

# UPM

## YM Raja Azilna Raja Mahmood, Mrs.

Faculty of computer science and IT technology University Putra Malaysia (Assessor)

> Abu Bakar Md. Sultan, PhD Dean School of Graduate Studies University Putra Malaysia Date:

## DECLARATION

## Declaration by a graduate student

I hereby confirm that:

- This thesis is my original work;
- Quotations, illustrations and citations have been duly referenced;
- This thesis has not been submitted previously or concurrently for any other Degree at any other institutions;
- Intellectual property from the thesis and copyright of thesis are fully-owned by University Putra Malaysia, as according to the University Putra Malaysia. (Research) Rules 2012;
- Written permission must be obtained from supervisor and the office of Deputy Vice-Chancellor (Research and Innovation) before thesis is published (in the form of written, printed or in electronic form) including books, journals, modules, proceedings, popular writings, seminar papers, manuscripts, posters, reports, lecture notes, learning modules or any other materials as stated in the University Putra Malaysia (Research) Rules 2012;
- There is no plagiarism or data falsification/fabrication in the thesis, and scholarly integrity is upheld as according to the University Putra Malaysia (Graduate Studies) Rules 2003(Revision 2012-2013) and the University Putra Malaysia (Research) Rules 2012. The thesis has undergone plagiarism detection software.

Signature:	 Date:
5	

Name and Matric No.: Awadh Salem Saleh Bajaher (GS48845)

#### ACKNOWLEDGEMENT

Firstly, all praises and gratitude go to the almighty Allah, whose favours are never ending for me, He is the one to give me all the abundances in my life to be become the person that I am today. After Allah, and what He has given me during my lifetime, I should not forget to mention those whom Allah put in my way to help me and support me. This thesis would not be completed without the help and supervision of my beloved supervisor **Dr. Nor Asila Wati Abdul Hamid,** who provided me with support, time, encouragement, unending guidance, and productive discussion.

I would like to thank my colleagues and the faculty team members, Department of Computer Science and the School of Graduate Studies staff at University Putra Malaysia. In the end, I would not forget to mention and thank eternally my father and mother who were the backbone for me and thank them for their encouragement, support, help and their spiritual and substantial support to continue and complete my master journey. No doubt that they deserve all my thanks because they are the reason for me to be here today. May Allah bless them the happiness, wellbeing and in the end the highest level of Jannah.

## DEDICATION

I am dedicating this thesis to my family and all my friends. Starting with the utmost gratitude to my parents who never failed to give me all types of support. My brother and my sisters who never forget to support me all along.

I likewise dedicate this thesis to many of my friends who always had been by my side whenever I needed them. I will not forget what was given to me by them either morally or physically, especially **Khalid Cawl** for empowering me to get a hand on my paper, and improve my work on CloudSim and **Daouda Camara** for helping me understand the algorithm details in spite of his busyness with his work.

Thanks to all of you

Thank you for being a huge part of this work

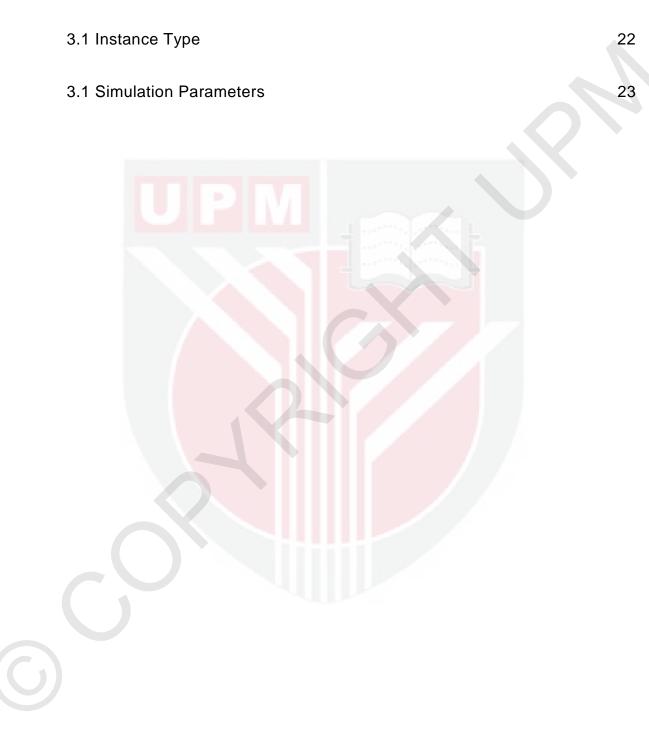
## List of Abbreviations

Virtual Machine
Cloudlet
Data Center Broker
Proportional Deadline Constrained
Deadline Constrained Critical Path
Service Provider
Million Instruction Per Second
Quality of Service
Virtual Machines
Million Instruction
Cloud Information Service
Data Center
Physical Machine
Earliest Completion Time

## List of Figures

3.1 The CloudSim Model	13.
3.2 Workflow Structure Example	16
3.3 The Sight Bioinformatics Workflow	17.
3.4 The Cybershake workflow	18
3.5 The Epigenomics Workflow	18
4.1 The PDC algorithm	28
5.1 Cybershake Workflow	31
5.2 Sight Workflow	32
5.3 Epigenomics Workflow	33

## List of Tables



СОР	YRIG	HTI
Abst	ract.	
APP	ROV	AL FORMIV
DEC		TION
ACK	Now	LEDGEMENT
List	of Ab	breviations
List	of Fig	guresIX
List	of Ta	blesX
Tabl	e of C	ContentsXI
1	INTF	RODUCTION
	1.1	Background1
	1.2	Problem Statement
	1.3	Objective of Research
	1.4	Scope of Research4
	1.5	Thesis Organization4
2	LITE	RATURE REVIEW5
	2.1	Introducation5

## **Table of Contents**

	2.2	Task Scheduling6
	2.3	Constrained Deadline Scheduling7
		2.3.1 Budget Constrained7
		2.3.2 Deadline Constrained
		2.3.3 Budget and Deadline Constrained8
3	мет	HODOLOGY11
	3.1	Introduction11
	3.2	The CloudSim Model
		3.2.1 Cloud Information Service (CIS)11
		3.2.2 Data center (DC)
		3.2.3 The Datacenter Broker (DCB)
		3.2.4 Virtual machine (VM)12
		3.2.5 The Cloudlet
		3.2.6 VM scheduler
		3.2.7 VM allocation13
		3.2.8 CloudSim model:
	3.3	The Problem Definition15
	3.4	Scientific Workflows
		3.4.1 Sight Workflow

		3.4.2 Cybershake Workflow17
		3.4.3 Epigenomics Workflow18
	3.5	Discrete Event Simulation (DES)19
	3.6	Datasets20
		3.6.1 The Parameters Used20
4	<b>IMP</b> I 4.1	LEMENTATION
		Implementation of The PDC Algorithm
		4.2.1 Workflow Leveling
		4.2.2 Deadline Distribution
		4.2.3 Task Selection
		4.2.4 Instance Selection
		4.2.5 The full PDC Algorithm
	4.3	The Performance Evaluation
5	RES	SULTS AND DISCUSSION
	5.1	Introduction
	5.2	Success Rate
		5.2.1 Cybershake
		5.2.2 Sight Workflow

5.2.3 Epigenomics Workflow
5.3 Normalized Cost (NC):
5.3.1 Cybershake workflow:
5.3.2 Sight Workflow:
5.3.3 Epigenomics Workflow:
CONCLUSION
References
6 Appendix A
7 Appendix B

## CHAPTER 1

## INTRODUCTION

## 1.1 Background

Cloud computing is the most modern trend, not only in computer science, but in a variety of sciences and fields. It is used for sharing resources, computations, deployments and experiments of different models. Moreover, cloud computing empowers noteworthy computational leverage to have it connected to numerous real-world issues, be they logical, industrial or therapeutic. When taking a scheduling viewpoint, it is concluded that the most attracting subset is the one related to multi-stage processing operations which might be denoted as workflow(Vahid Arabnejad, Bubendorfer, & Ng, 2017).

More so, commercial clouds are also utilized for analysis of scientific works and computations due to its scalability, latest high-quality hardware as well as pay-as-you-go cost model. Commercial clouds are based on a pay-peruse model. Most of the services of commercial clouds are paid based on duration of use of their resources like storage and network bandwidth. Furthermore, scientific computations have been lately performed on the cloud using the commercial cloud environments such as Amazon EC2, Google. etc. Such scientific computations usually involve fault tolerance, load balancing and accessing certain resources like GPU. Nevertheless, such a flexible cloud model could cause high costs whenever insufficient scheduling are performed (Chard et al., 2015).

This work tends to cover the issue of scheduling deadline-constrained for scientific workflows such as Cybershake, Sight and Epigenomics which are discussed in details in the methodology. Additionally, a number of novel scheduling deadline constrained algorithms were introduced.

The Proportional Deadline Constrained (PDC) algorithm has been presented by Arabnejad, Bubendofer, and Ng to solve the issue of scheduling deadline constrained schedules in commercial clouds which usually results in high failure rate as well as high cost (2017). Besides, the PDC algorithm has four main stages, namely workflow levelling, deadline distribution, task selection and instance selection. Those four stages are elaborated more in the methodology. After implementing all four stages, the PDC algorithm will then be compared to existing algorithms such as IC-PCP. PDC is expected to show higher success rate.

The rest of the work consists of subsections, problem statement, objectives and scope of the work. While the main sections, the literature review, methodology, implementation, results and analysis and conclusion.

2

## 1.2 **Problem Statement**

Commercial clouds can be easily accessed globally due to the rapid development of technology. Furthermore, commercial cloud corporations usually charge their clients on hourly basis as their resources are being used such as network bandwidth, storage, etc. It is also known as pay-per-use model. Moreover, executing scientific workflows on commercial clouds using pay-as-you-go model requires certain algorithms to schedule and complete tasks to pay as less as possible. There have been several studies presenting new algorithms to generate deadline constrained schedules to minimize the execution cost as well as the high failure rate in schedule constructions. However, there are increased failure rates whenever tight deadlines are produced (Vahid Arabnejad et al., 2017).

## 1.3 Objective of Research

The main objective of this project is the following:

- To evaluate the performance of Proportional Deadline Constrained (PDC) algorithm in scheduling scientific workflow.

#### 1.4 Scope of Research

The focus of the study is the issue of scheduling tight-deadline scientific workflow, studying how the resources are provisioned and utilized for executing certain tasks parallelly. Also, the evaluation of the effectiveness of PDC algorithm in scheduling scientific workflow. Additionally, the evaluation will be performed using discrete event simulation represented in CloudSim framework for modelling and cloud computing simulation.

## 1.5 Thesis Organization

The thesis is organized as the following:

**Chapter 1** Introduction that elaborates the cloud computing environment, commercial clouds and the summary of work.

**Chapter 2** The literature review focuses on explaining more about the cloud and commercial clouds. Also, it includes a critically reviewed related work to this study.

**Chapter 3** The methodology illustrated more about the definition of the problem, the algorithm, the method used, the parameters and workloads included as well as the simulator used.

**Chapter 4** Implementation which shows how algorithms works, how this study simulation is performed and the experimentation executed.

**Chapter 5** Results and Discussion illustrates the results of the simulation done in Chapter 4 as well as analyzing it.

**In chapter 6** The conclusion includes a summary of the whole work.

#### References

- Arabnejad, V., & Bubendorfer, K. (2015). Cost Effective and Deadline Constrained Scientific Workflow Scheduling for Commercial Clouds. In 2015 IEEE 14th International Symposium on Network Computing and Applications (pp. 106– 113). https://doi.org/10.1109/NCA.2015.33
- Arabnejad, Vahid, Bubendorfer, K., & Ng, B. (2017). Scheduling deadline constrained scientific workflows on dynamically provisioned cloud resources. *Future Generation Computer Systems*, 75, 348–364. https://doi.org/10.1016/j.future.2017.01.002
- Belalem, G., Tayeb, F. Z., & Zaoui, W. (2010). Approaches to Improve the Resources Management in the Simulator CloudSim. In Information Computing and Applications (pp. 189–196). Springer, Berlin, Heidelberg. https://doi.org/10.1007/978-3-642-16167-4\_25
- Buyya, R., Yeo, C. S., Venugopal, S., Broberg, J., & Brandic, I. (2009). Cloud computing and emerging IT platforms: Vision, hype, and reality for delivering computing as the 5th utility. *Future Generation Computer Systems*, 25(6), 599–616. https://doi.org/10.1016/j.future.2008.12.001

Calheiros, R. N., Ranjan, R., Beloglazov, A., De Rose, C. A. F., & Buyya, R. (2011).
CloudSim: a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms. *Software: Practice and Experience*, *41*(1), 23–50. https://doi.org/10.1002/spe.995

- Calheiros, R. N., Ranjan, R., Beloglazov, A., & Rose, A. F. De. (2011). CloudSim : a toolkit for modeling and simulation of cloud computing environments and evaluation of resource provisioning algorithms, (August 2010), 23–50. https://doi.org/10.1002/spe
- Chard, R., Chard, K., Bubendorfer, K., Lacinski, L., Madduri, R., & Foster, I. (2015). Cost-Aware Cloud Provisioning (pp. 136–144). IEEE. https://doi.org/10.1109/eScience.2015.67
- Choi, B. K., & Kang, D. (2013). *Modeling and Simulation of Discrete Event Systems* (1 edition). Hoboken, NJ: Wiley.
- Gulati, A., & Chopra, R. K. (2013). Dynamic Round Robin for Load Balancing in a Cloud Computing, (6), 5.
- Haidri, R. A., Katti, C. P., & Saxena, P. C. (2017). Cost effective deadline aware scheduling strategy for workflow applications on virtual machines in cloud computing. *Journal of King Saud University - Computer and Information Sciences*. https://doi.org/10.1016/j.jksuci.2017.10.009
- Juve, G., Chervenak, A., Deelman, E., Bharathi, S., Mehta, G., & Vahi, K. (2013).
  Characterizing and profiling scientific workflows. *Future Generation Computer* Systems, 29(3), 682–692.
  https://doi.org/10.1016/j.future.2012.08.015
- Malawski, M., Figiela, K., Bubak, M., Deelman, E., & Nabrzyski, J. (2015).
  Scheduling Multilevel Deadline-constrained Scientific Workflows on Clouds
  Based on Cost Optimization. Sci. Program., 2015, 5:5–5:5.
  https://doi.org/10.1155/2015/680271

Margaret Rouse. (n.d.). What is cloud computing? - Definition from WhatIs.com.RetrievedJuly5,2018,fromhttps://searchcloudcomputing.techtarget.com/definition/cloud-computing

Masoudi Khorsand, M., Jamali, S., & Analoui, M. (2014). A Survey of Job Scheduling Algorithms Whit Hierarchical Structure to LoadBalancing In Grid Computing Environments. *International Journal of Computer Applications Technology and Research*, 3, 68–72. https://doi.org/10.7753/IJCATR0301.1015

- Rodriguez, M. A., & Buyya, R. (2014). Deadline Based Resource Provisioningand
   Scheduling Algorithm for Scientific Workflows on Clouds. *IEEE Transactions on Cloud Computing*, 2(2), 222–235.
   https://doi.org/10.1109/TCC.2014.2314655
- Roy, S., Banerjee, S., Chowdhury, K. R., & Biswas, U. (2017). Development and analysis of a three phase cloudlet allocation algorithm. *Journal of King Saud University Computer and Information Sciences*, 29(4), 473–483. https://doi.org/10.1016/j.jksuci.2016.01.003
- Sakellariou, R., Zhao, H., Tsiakkouri, E., & Dikaiakos, M. D. (2007). Scheduling
  Workflows with Budget Constraints. In *Integrated Research in GRID Computing* (pp. 189–202). Springer, Boston, MA.
  https://doi.org/10.1007/978-0-387-47658-2\_14
- Ullman, J. D. (1975). NP-complete scheduling problems. Journal of Computer and System Sciences, 10(3), 384–393. https://doi.org/10.1016/S0022-0000(75)80008-0

- Wu, F., Wu, Q., & Tan, Y. (2015). Workflow scheduling in cloud: a survey. The Journal of Supercomputing, 71(9), 3373–3418. https://doi.org/10.1007/s11227-015-1438-4
- Wu, Z., Liu, X., Ni, Z., Yuan, D., & Yang, Y. (2013). A market-oriented hierarchical scheduling strategy in cloud workflow systems. *The Journal of Supercomputing*, 63(1), 256–293. https://doi.org/10.1007/s11227-011-0578-4
- Yuan, Y., Li, X., Wang, Q., & Zhu, X. (2009). Deadline Division-based Heuristic for Cost Optimization in Workflow Scheduling. *Inf. Sci.*, 179(15), 2562– 2575. https://doi.org/10.1016/j.ins.2009.01.035
- Zheng, W., & Sakellariou, R. (2013). Budget-Deadline Constrained Workflow Planning for Admission Control. *Journal of Grid Computing*, 11(4), 633–651. https://doi.org/10.1007/s10723-013-9257-4