

**A BIG DATA ANALYTIC FRAMEWORK OVER  
FEDERATED DATA CENTERS FOR INTELLIGENT  
TRAVEL RECOMMENDERS**

H.P.I. Udayanthi

(188067X)

Degree of Master of Science

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

February 2021

**A BIG DATA ANALYTIC FRAMEWORK OVER  
FEDERATED DATA CENTERS FOR INTELLIGENT  
TRAVEL RECOMMENDERS**

Hewa Pathirange Ishara Udayanthi

(188067X)

Thesis submitted in partial fulfilment of the requirements for the degree Master of  
Science in Computational Mathematics

Department of Computational Mathematics

University of Moratuwa

Sri Lanka

February 2021

## **Declaration**

I declare that this is my own work and this thesis does not incorporate without acknowledgement any material previously submitted for a Degree or Diploma in any other University or institute of higher learning, and to the best of my knowledge and belief, it does not contain any material previously published or written by another person except where the acknowledgement made is made in the text.

Also, I hereby grant to the University of Moratuwa the non-exclusive right to reproduce and distribute my thesis, in whole or part in print, electronic or another medium. I retain the right to use this content in whole or part in future works (such as articles or books)

Signature:

Date:

The above candidate has carried out the research for the Master's thesis (By research) under my supervision.

Signature of the supervisor:

Date:

**Dedicated to**

*My beloved Father and Mother*

## **Acknowledgements**

Many people have helped their best to successfully complete this research. I acknowledge all of them for their valuable thoughts and constant encouragement gave me to make my project a reality.

I would first like to acknowledge my supervisor, Senior Lecturer Thushari Silva, for accepting me as his research student and giving excellent support and advice. Dr. Thushari Silva is a great mentor whose expertise was invaluable in formulating the research questions and methodology. Your insightful feedback pushed me to sharpen my thinking and brought my work to a higher level.

Exceptional thanks goes to Dr. A. Shehan Perera for his valuable comments and guidance as my examiner of progress review panels.

I wish to extend my sincere thanks for all the academic and non-academic staff of the Department of Computational Mathematics, the University of Moratuwa for their kind-hearted help to fulfil my research work.

I acknowledge the sacrificial dedication of my family members for their encouragement and cooperation by managing all the works while I was busy with my research work.

## Abstract

Big Data is a series of enormous and complex data sets that are nearly impossible to store and process using traditional data storing and processing methods. The emergence of heterogeneous data in different domains causes significant challenges in data manipulation and decision making. In recent years, the requirement for analysis of heterogeneous data on distributed data storages has been increased and has gained a lot of researchers' attention. Distributed data storage systems and parallel data processing techniques are typically used for data-intensive computing today. Due to the rapid growth of data, a single-cluster environment is inadequate to process that much data. At the same time, there are heterogeneous data sources on different platforms, which need to inter-connect to derive meaningful analysis. The MapReduce software paradigm has surfaced to fill the gap, and it has been successfully operating on systems. However, only single cluster environments are supported by the current implementation of MapReduce and cannot be applied to federated heterogeneous data centers. Hence, it does not have enough capabilities to process heterogeneous data sources. This research presents a big data analytic framework that supports the integration of heterogeneous data sources on distributed computing models across different data centers. The architecture of this framework is based on a master/slave distributed computing model and Map - Reduce - Merge - GlobalReduce is presented as the programming model. Besides, the performance of the novel approach is measured under different cluster configurations, and experimental evaluations had shown promising results for the proposed framework compared to a single cluster environment.

**Keywords:** Big Data, Heterogeneous data centers, Hierarchical MapReduce, Cloud computing, Multi-cluster

# TABLE OF CONTENTS

<b>CHAPTER 1 - INTRODUCTION</b>	<b>1</b>
1.1 Prolegomena	1
1.2 Aims and Objectives	2
1.3 Background and Motivation	2
1.4 Research Statement	7
1.5 Current Approaches in Developing Hierarchical Big Data Analytic Framework	8
1.6 The Proposed Solution	11
1.7 Resource requirements	12
1.8 Organization	12
1.9 Summary	13
<b>CHAPTER 2 - LITERATURE REVIEW</b>	<b>14</b>
2.1 Introduction	14
2.2 Hadoop Architecture	14
2.2.1 Hadoop Distributed File System (HDFS)	15
2.2.2 Hadoop MapReduce	17
2.2.3 Hadoop YARN	18
2.3 Hierarchical MapReduce Frameworks	18
2.3.1 HDC-Hadoop	19
2.3.2 G-Hadoop	21
2.3.3 Federated MapReduce (Fed-MR)	22
2.3.4 Map-Reduce-GlobalReduce	24
2.3.5 Map-Reduce-Merge	26
2.4 Summary	27
<b>CHAPTER 3 - RESEARCH METHODOLOGY</b>	<b>29</b>
3.1 Introduction	29
3.2 Architecture and Execution Flow	29
3.3 Programming Model	30

3.4	Workflow	36
3.5	Summary	37
	<b>CHAPTER 4 – CASE STUDY</b>	<b>38</b>
4.1	Introduction	38
4.2	Recommender Systems	38
4.3	Popularity Based Intelligent Travel Recommendation System	41
4.4	Data Sources	41
4.4.1	Tweets	42
4.4.2	Ratings	42
4.5	Summary	42
	<b>CHAPTER 5 – PERFORMANCE EVALUATION</b>	<b>43</b>
5.1	Introduction	43
5.2	Experimental Setup	43
5.3	Implementation	44
5.4	Performance Comparison	45
5.5	Summary	48
	<b>CHAPTER 6 – CONCLUSION AND FUTURE WORKS</b>	<b>49</b>
6.1	Introduction	49
6.2	Developing Big Data Analytic Framework over federated data centers	49
6.3	Objectives-wise Achievement	50
6.4	Limitations and Future Work	51
6.5	Summary	52
	<b>REFERENCES</b>	<b>53</b>
	<b>APPENDIX A - DATA SETS</b>	<b>58</b>
	<b>APPENDIX B - IMPLEMENTATION OF THE FRAMEWORK</b>	<b>59</b>





## LIST OF FIGURES

Figure 1.1	Characteristic of Big Data .....	3
Figure 2.1	Hadoop Framework .....	14
Figure 2.2	HDFS Architecture .....	17
Figure 2.3	Architecture overview of HDC-Hadoop.....	20
Figure 2.6	Basic Architecture of G-Hadoop .....	22
Figure 2.4	Federated MapReduce operational model .....	23
Figure 2.5	Map-Reduce - Global Reduce Framework Architecture .....	26
Figure 2.7	Data and Control Flow of Map - Reduce – Merge .....	27
Figure 3.1	High-level Architecture .....	30
Figure 3.2	Data Flow Sequence .....	32
Figure 3.3	Map function for tweets algorithm .....	34
Figure 3.4	Reduce function for tweets algorithm.....	34
Figure 3.5	Map function for ratings algorithm .....	35
Figure 3.6	Reduce function for ratings algorithm.....	35
Figure 3.7	Merge algorithm .....	35
Figure 3.8	Global Reduce algorithm.....	36
Figure 3.9	Workflow .....	37
Figure 4.1	Execution Flow of Recommendation System .....	41
Figure 5.1	Execution Flow and Results .....	44
Figure 5.2	Total Execution Time under Different Cluster Configurations.....	45
Figure 5.3	Performance Breakdown .....	46

## LIST OF TABLES

Table 3.1 Input and output formats of Map, Reduce, Merge, and Global Reduce functions.....	31
Table 5.1 Precision, Recall and F-measure .....	47

## **LIST OF ABBREVIATIONS**

HDFS – Hadoop Distributed File System

YARN – Yet Another Resource Negotiator

GPS – Global Positioning System

RDMS – Relational Database Management System

RAM – Random Access Memory

OS – Operating System

CPU – Central Processing Unit

i.e., That is

NLTK – Natural Language Toolkit

NLP – Natural Language Processing

LAN – Local Area Network

URL – Uniform Resource Locator