# CLOUD-BASED DOCUMENT MANAGEMENT FRAMEWORK FOR BUSINESS PROCESS OPTIMIZATION

A. Nalaka Arjuna Premathilaka

(158240U)

Degree of Master of Science in Computer Science

Department of Computer Science and Engineering

University of Moratuwa
Sri Lanka

February 2019

# CLOUD-BASED DOCUMENT MANAGEMENT FRAMEWORK FOR BUSINESS PROCESS OPTIMIZATION

A. Nalaka Arjuna Premathilaka

(158240U)

Thesis submitted in partial fulfillment of the requirements for the degree

Master of Science in Computer Science

Department of Computer Science and Engineering

University of Moratuwa

Sri Lanka

February 2019

# DECLARATION

I declare that this is my own work and this dissertation does not incorporate without acknowledgment 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 acknowledgment is made in the text.

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

Signature: ……………………. Date: ……………….

Name: A. Nalaka Arjuna Premathilaka (158240U)

The supervisor/s should certify the dissertation with the following declaration.

The above candidate has carried out research for the Master's Dissertation under my supervision.

Signature of the supervisor: …………………………. Date: ………………..

Name: Dr. Indika Perera

# ABSTRACT

Nowadays, most of the corporates and private users are tend to use the Cloud services because of its benefits such as cost reduction, flexibility & scalability, high availability, accessibility and many more. When organizations upload their sensitive data to the public cloud storage in plain text, the main concern is data security & privacy. Cloud data must be protected from the external attackers, intruders and from Cloud storage owners as well. There are few examples that even the cloud storage providers were involved in the security breaches of the data in their storages. Therefore, In order to achieve the data security and privacy in the public cloud storages, usually data will be encrypted prior upload to the cloud. However, in public cloud storages such as Dropbox, Amazon S3, Mozy, and others, perform data deduplication to save space by removing the repeated chunks of the files or data blocks. This will help to reduce storage usage and has a cost benefit as well. But when the data is encrypted, the de-duplication process is not working as expected and storage space savings are lost.

This research focuses on identifying a method to overcome these issues in public cloud storages when storing sensitive data. This research implements a solution/framework for corporate and individual users to use public cloud storage by ensuring data security and space saving as well. Implementation of this research will introduce an additional application layer between public cloud storage and cloud users. It handles communication between the user and the cloud storage and performs the data de-duplication and encryption prior to uploading data chunks to the public cloud.

**Keywords**: Cloud storage, Data security, Data de-duplication, Data encryption, Public cloud, Cryptography, Access control, Data Privacy, Authorization

# ACKNOWLEDGMENTS

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| | |
|---|---|
| QOS | Quality of Service |
| GB | Giga Byte |
| KB | Kilo Byte |
| Db | Database |
| IaaS | Infrastructure as a Service |
| PaaS | Platform as a Service |
| SaaS | Software as a Service |
| CE | Convergent Encryption |
| AES | Advanced Encryption System |
| CBC | Cipher Block Chaining |
| ECB | Electronic Code Book |