



# Reliability Assurance of Big Data in the Cloud

Cost-Effective Replication-Based Storage

Wenhao Li, Yun Yang, Dong Yuan

# Reliability Assurance of Big Data in the Cloud

Cost-Effective Replication-Based  
Storage

*by*

***Wenhao Li, Yun Yang, Dong Yuan***

**Centre for Computing and Engineering**

**Software Systems,**

**School of Software and Electrical Engineering**

**Swinburne University of Technology**

**Hawthorn, Melbourne, Australia**



**ELSEVIER**

AMSTERDAM • BOSTON • HEIDELBERG • LONDON  
NEW YORK OXFORD • PARIS • SAN DIEGO  
SAN FRANCISCO • SINGAPORE SYDNEY • TOKYO

Morgan Kaufmann is an imprint of Elsevier



Acquiring Editor: Todd Green  
Editorial Project Manager: Lindsay Lawrence  
Project Manager: Surya Narayanan Jayachandran  
Designer: Mark Rogers

Morgan Kaufmann is an imprint of Elsevier  
225 Wyman Street, Waltham, MA 02451, USA

Copyright © 2015 Elsevier Inc. All rights reserved.

No part of this publication may be reproduced or transmitted in any form or by any means, electronic or mechanical, including photocopying, recording, or any information storage and retrieval system, without permission in writing from the publisher. Details on how to seek permission, further information about the Publisher's permissions policies and our arrangements with organizations such as the Copyright Clearance Center and the Copyright Licensing Agency, can be found at our website: [www.elsevier.com/permissions](http://www.elsevier.com/permissions).

This book and the individual contributions contained in it are protected under copyright by the Publisher (other than as may be noted herein).

### Notices

Knowledge and best practice in this field are constantly changing. As new research and experience broaden our understanding, changes in research methods, professional practices, or medical treatment may become necessary.

Practitioners and researchers must always rely on their own experience and knowledge in evaluating and using any information, methods, compounds, or experiments described herein. In using such information or methods they should be mindful of their own safety and the safety of others, including parties for whom they have a professional responsibility.

To the fullest extent of the law, neither the Publisher nor the authors, contributors, or editors, assume any liability for any injury and/or damage to persons or property as a matter of products liability, negligence or otherwise, or from any use or operation of any methods, products, instructions, or ideas contained in the material herein.

### British Library Cataloguing-in-Publication Data

A catalogue record for this book is available from the British Library

### Library of Congress Cataloging-in-Publication Data

A catalog record for this book is available from the Library of Congress

ISBN: 978-0-12-802572-7

For information on all MK publications  
visit our website at [www.mkp.com](http://www.mkp.com)

This book has been manufactured using Print On Demand technology. Each copy is produced to order and is limited to black ink. The online version of this book will show color figures where appropriate.



Working together  
to grow libraries in  
developing countries

[www.elsevier.com](http://www.elsevier.com) • [www.bookaid.org](http://www.bookaid.org)

# **Reliability Assurance of Big Data in the Cloud**

**Cost-Effective Replication-Based  
Storage**

# About the Authors

**Wenhao Li** received the BEng and MEng degrees from Shandong University, China, in 2007 and 2010, respectively, and PhD degree from the School of Software and Electrical Engineering, Swinburne University of Technology, Melbourne, Australia, in 2014, supervised by Professor Yun Yang. All his degrees are held in computer science. He is currently a SSEE research fellow in the School of Software and Electrical Engineering at Swinburne University of Technology, Melbourne, Australia. His research interests include parallel and distributed computing, cloud and grid computing, work flow technologies, and data management in distributed computing environments. He is a member of the Institute of Electrical and Electronics Engineers (IEEE).



**Yun Yang** received the BS degree from Anhui University, China, in 1984, the MEng degree from the University of Science and Technology of China, in 1987, and the PhD degree from the University of Queensland, Australia, in 1992, all in computer science. He is currently a full professor in the School of Software and Electrical Engineering at Swinburne University of Technology, Melbourne, Australia. Prior to joining Swinburne in 1999 as an associate professor, he was a lecturer and senior lecturer at Deakin University, Australia, during 1996–1999. From 1993 to 1996 he was a (senior) research scientist at DSTC Cooperative Research Centre for Distributed Systems Technology, Australia. He was also at Beijing University of Aeronautics and Astronautics, China, during 1987–1988. He has coauthored four books, published more than 200 papers in journals, and refereed conference proceedings. He is currently on the Editorial Board of IEEE Transactions on Cloud Computing. His current research interests include software technologies, cloud computing, p2p/grid/cloud workflow systems, and service-oriented computing. He is a senior member of the IEEE.



**Dong Yuan** received the BEng and MEng degrees from Shandong University, Jinan, China, in 2005 and 2008, respectively, and the PhD degree from Swinburne University of Technology, Melbourne, Australia, in 2012, all in computer science. He is currently a research fellow in SSEE to the School of Software and Electrical Engineering at Swinburne University of Technology, Melbourne, Australia. His research interests include data management in parallel and distributed systems, scheduling and resource management, and grid and cloud computing.



# Preface

Cloud computing is the latest distributed computing paradigm that provides redundant, inexpensive, and scalable resources in a pay-as-you-go fashion to meet various application requirements. Nowadays, with the rapid growth of Cloud computing, the size of Cloud data is expanding at a dramatic speed. A huge amount of data that are big in sizes and large in amounts are generated and processed by Cloud applications with data-intensive characteristics. For maintaining big data in the cloud, data reliability-related issues are considered more important than ever before. However, current data storage and data reliability-ensuring strategies based on multiple replicas have become a bottleneck for big data storage in the Cloud. For storing massive data in the Cloud, such strategies could consume a huge amount of storage resources on replication, and hence incur a huge storage cost and cause negative effects for both Cloud storage providers and storage users. Therefore, a higher demand has been put forward to Cloud storage. While the requirement of data reliability should be met in the first place, data in the Cloud needs to be stored in a highly cost-effective manner.

In this book, we investigate the trade-off between data storage cost and data reliability assurance for big data in the Cloud. The research is motivated by a scientific application for astrophysics pulsar searching surveys, which is of typical data-intensive characteristics and contains complex and time-consuming tasks that process hundreds of terabytes of data. To store the massive amount of application data into the Cloud, our novel research stands from the Cloud storage service providers' perspective and investigates the issue on how to provide cost-effective data storage while meeting the data reliability requirement throughout the whole Cloud data life cycle. Our research in this book presents four major contributions. According to different stages within the Cloud data life cycle, these four contributions are presented in the following sequence.

1. For describing data reliability in the Cloud, a novel generic data reliability model is proposed. Based on a Cloud with replication-based data storage scheme, the data reliability model is able to describe the reliability of the Cloud data throughout their life cycles, in which they are stored with different redundancy levels and stored on different storage devices in different stages respectively. Compared with existing data reliability models that assume a constant disk failure rate, our generic data reliability model is able to better describe data reliability over a wide range of failure rate patterns of storage devices.
2. To facilitate data creation, a minimum replication calculation approach for meeting a given data reliability requirement is proposed. Based on the data reliability model, this approach calculates the minimum number of replicas that needs to be created for meeting certain data reliability requirement and predicts the reliability of the data stored for a certain amount of time. In addition, the minimum replication can also act as a benchmark, which can be used for evaluating the cost-effectiveness of various replication-based data storage approaches.
3. In the data maintenance stage, in order to maintain the Cloud data with the minimum replication level, a novel cost-effective data reliability assurance mechanism named Proactive

Replica Checking for Reliability (PRCR) is proposed. Based on the minimum replication that is created, PRCR is able to maintain the huge amounts of Cloud data with negligible overhead, while a wide variety of data reliability assurance can be provided. Compared with the widely used conventional three-replica data storage and data reliability-ensuring strategy, PRCR significantly lowers storage cost in the Cloud. PRCR can reduce from one-third to two-thirds of the Cloud storage space consumption. Even more saving can be achieved compared with data storage strategies with higher replication levels.

4. In the data creation and recovery stages, in order to reduce the data transfer cost, a cost-effective strategy named Link Rate Controlled Data Transfer (LRCDT) is proposed. By scheduling bandwidth in a link rate controlled fashion, LRCDT could significantly reduce the energy consumption during the data creation and recovery process in the Cloud network. The result in our simulation indicates that LRCDT is able to reduce energy consumption by up to 63% when compared to existing data transfer strategies.

The research issue of this book is significant and has practical value to the Cloud computing technology. Especially, for Cloud applications that are of data-intensive characteristics, our research could significantly reduce their storage cost while meeting the data reliability requirement and have a positive effect on promoting the development of Cloud.

# Acknowledgments

This research is partly supported by the Australian Research Council (ARC) grants under DP110101340 and LP130100324. We are grateful for English proofreading by Z. Sterling.



# Contents

<b>About the Authors</b>	<b>ix</b>
<b>Preface</b>	<b>xi</b>
<b>Acknowledgments</b>	<b>xiii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Data reliability in the Cloud	1
1.2 Background of Cloud storage	2
1.2.1 Distinctive features of Cloud storage systems	2
1.2.2 The Cloud data life cycle	4
1.3 Key issues of research	5
1.4 Book overview	6
<b>2 Literature review</b>	<b>9</b>
2.1 Data reliability assurance in hardware	9
2.1.1 Disk	9
2.1.2 Other storage medias	12
2.2 Data reliability assurance in software	13
2.2.1 Replication for data reliability	13
2.2.2 Erasure coding for data reliability	14
2.3 Data transfer for distributed systems	15
2.4 Summary	17
<b>3 Motivating example and problem analysis</b>	<b>19</b>
3.1 Motivating example	19
3.1.1 The pulsar searching application process	19
3.1.2 The pulsar searching application data flow	21
3.1.3 Storing pulsar searching data in the Cloud	23
3.2 Problem analysis	24
3.2.1 Two major factors of Cloud storage cost	24
3.2.2 Data storage devices and schemes	25
3.2.3 Cloud network and data transfer activities	26
3.2.4 Research issues	27
3.3 Summary	29
<b>4 Generic data reliability model in the cloud</b>	<b>31</b>
4.1 Properties of the data reliability model	31
4.1.1 Reliability metrics	31
4.1.2 Data reliability model type	32
4.1.3 Failure rate pattern of storage devices	33

4.2	Generic data reliability model	33
4.2.1	Data reliability with static disk failure rate	33
4.2.2	Data reliability with variable disk failure rate	34
4.2.3	Generic data reliability model for multi-replicas	35
4.3	Summary	36
<b>5</b>	<b>Minimum replication for meeting the data reliability requirement</b>	<b>37</b>
5.1	The minimum replication calculation approach	37
5.1.1	Minimum replication calculation formulas	37
5.1.2	Optimization of the minimum replication calculation formulas	40
5.2	Minimum replication benchmark	41
5.3	Evaluation of the minimum replication calculation approach	42
5.4	Summary	43
<b>6</b>	<b>Cost-effective data reliability assurance for data maintenance</b>	<b>45</b>
6.1	Proactive replica checking	45
6.2	Overview of PRCR	46
6.2.1	User interface	47
6.2.2	PRCR node	48
6.3	Working process of PRCR	49
6.4	Optimization algorithms in PRCR	50
6.4.1	Minimum replication algorithm	50
6.4.2	Metadata distribution algorithm	52
6.5	Evaluation of PRCR	54
6.5.1	Performance of PRCR	55
6.5.2	Cost-effectiveness of PRCR	59
6.5.3	Summary of the evaluation	62
6.6	Summary	62
<b>7</b>	<b>Cost-effective data transfer for data creation and data recovery</b>	<b>65</b>
7.1	Determining the deadline for data creation and data recovery	65
7.2	Cloud network model	66
7.2.1	Overall network model	67
7.2.2	Pipeline model	67
7.2.3	Pipeline agenda model	68
7.2.4	Overall agenda model	68
7.3	Energy consumption model for Cloud data transfer	69
7.4	Novel cost-effective data transfer strategy LRCDT	70
7.5	Evaluation of LRCDT	74
7.5.1	Parameters of simulation	74
7.5.2	Energy consumption comparison	75
7.5.3	Task completion time comparison	77
7.6	Summary	78

---

<b>8</b>	<b>Conclusions and future work</b>	<b>79</b>
8.1	Summary of this book	79
8.2	Key contributions of this book	80
8.3	Further discussion and future work	81
8.3.1	Further discussions	81
8.3.2	Future work	82
	<b>Bibliography</b>	<b>83</b>
	<b>Appendix</b>	<b>87</b>
	<b>Index</b>	<b>89</b>

With the rapid growth in the size of Cloud data, cost-effective data storage has become one of the key issues in Cloud research, yet the reliability of the huge amounts of Cloud data needs to be fully assured. In this book, we investigate the trade-off of cost-effective data storage and data reliability assurance in the Cloud. The novel research stands from the Cloud storage service providers' perspective and investigates the issue on how to provide cost-effective data storage service while meeting the data reliability requirement throughout the whole Cloud data life cycle. This topic is important and has a practical value to Cloud computing technology. Especially, for data-intensive applications that are of data-intensive characteristics, our research could dramatically reduce its storage cost while meeting the data reliability requirement and hence has a positive impact on promoting the deployment of the Cloud.

This chapter introduces the background knowledge and key issues of this research. It is organized as follows. Section 1.1 gives the definition of data reliability and briefly introduces current data reliability assurance technologies in the Cloud. Section 1.2 introduces the background knowledge related to Cloud storage. Section 1.3 outlines the key issues of the research. Finally, Section 1.4 presents an overview for the book structure.

## 1.1 Data reliability in the Cloud

The term “reliability” is widely used as an aspect of the service quality provided by hardware, systems, Web services, etc. In Standard TL9000, it is defined as “the ability of an item to perform a required function under stated conditions for a stated time period” [1]. For data reliability specifically, which refers to the reliability provided by the data storage services and systems for the stored data, it can be defined as “the probability of the data surviving in the system for a given period of time” [2]. While the term “data reliability” is sometimes used in the industry as a superset of data availability and various other topics, in this book we will stick to the definition of data reliability given earlier.

Data reliability indicates the ability of the storage system to keep data consistent, hence it is always one of the key metrics of a data storage/management system. In large-scale distributed systems, due to the big quantity of storage devices being used, failures of storage devices occur frequently [3]. Therefore, the importance of data reliability is prominent, and these systems need better design and management to cope with frequent failures. Increasing the data redundancy level could be a good way for increasing data reliability [4,5]. Among several major approaches for increasing the data redundancy level, data replication is currently the most popular approach in distributed storage systems. At present, data replication has been widely adopted in many

current distributed data storage/management systems in both industry and academia, which include examples such as OceanStore [6], DataGrid [7], Hadoop Distributed File System [8], Google File System [9], Amazon S3 [10], and so forth. In these storage systems, several replicas are created for each piece of data. These replicas are stored in different storage devices, so that the data have better chance to survive when storage device failures occur.

In recent years, Cloud computing is emerging as the latest distributed computing paradigm, which provides redundant, inexpensive, and scalable resources in a pay-as-you-go fashion to meet various application requirements [11]. Since the advent of Cloud computing in late 2007 [12], it has fast become one of the most promising distributed solutions in both industry and academia. Nowadays, with the rapid growth of Cloud computing, the size of Cloud storage is expanding at a dramatic speed. It is estimated that by 2015 the data stored in the Cloud will reach 0.8 ZB (i.e.,  $0.8 \times 10^{21}$  bytes or 800,000,000 TB), while more data are “touched” by the Cloud within their life cycles [13]. For maintaining such a large amount of Cloud data, data reliability in the Cloud is considered more important than ever before. However, due to the accelerating growth of Cloud data, current replication-based data reliability management has become a bottleneck for the development of Cloud data storage. For example, storage systems such as Amazon S3, Google File System, and Hadoop Distributed File System all adopt similar data replication strategies called the “conventional multi-replica replication strategy,” in which a fixed number of replicas (normally three) are stored for all data to ensure the reliability requirement. For storage of the huge amounts of Cloud data, these conventional multi-replica replication strategies consume a lot of storage resources for additional replicas. This could cause negative effects for both the Cloud storage providers and users. On one hand, from the Cloud storage provider’s perspective, the excessive consumption of storage resources leads to a big storage overhead and increases the cost for providing the storage service. On the other hand, from the Cloud storage user’s perspective, according to the pay-as-you-go pricing model, the excessive storage resource usage will finally be paid by the storage users. For data-intensive Cloud applications specifically, the incurred excessive storage cost could be huge. Therefore, Cloud-based applications have put forward a higher demand for cost-effective management of Cloud storage. While the requirement of data reliability should be met in the first place, data in the Cloud needs to be stored in a highly cost-effective manner.

## 1.2 Background of Cloud storage

In this section, we briefly introduce the background knowledge of Cloud storage. First, we introduce the distinctive features of Cloud storage systems. Second, we introduce the Cloud data life cycle.

### 1.2.1 Distinctive features of Cloud storage systems

Data reliability is closely related to the structure of the storage system and how the storage system is being used. Different from other distributed storage systems, the

Cloud storage system has some distinctive features that could either be advantages or challenges for the data reliability management of Cloud data.

#### *1.2.1.1 On-demand self-service and pay-as-you-go pricing model*

The on-demand usage of Cloud storage service and pay-as-you-go payment fashion have greatly facilitated the storage users that they only need to pay for the resources used for storing their data for a needed time period. The cost is easy to be estimated according to the size of data generated [14]. However, based on the pay-as-you-go model, every usage of the resources can be strictly reflected onto the bills payable at the end of the month. Therefore, minimizing resource consumption becomes demanding and critical. This principle is not only applicable to the service users, but also to the Cloud storage service providers. In most current Cloud storage services, excessive data redundancy is compulsorily generated to ensure data reliability. For data-intensive applications, such excessive data redundancy consumes a large amount of storage resources, and hence incurs a very high cost.

#### *1.2.1.2 Redundant and scalable virtualized resources*

In the Cloud, large amounts of virtualized computing and storage resources are pooled to serve users with various demands [1]. Redundant computing resources make it easy to conduct parallel processing, while the redundant storage resources make it easy to distribute data. For meeting a higher computing/storage demand, the resource pool can be scaled out rapidly, and the virtualization keeps the complex procedures transparent from the service users. However, the virtualization of resources has also led to a challenge that various kinds of data reliability requirement need to be fully assured to make the Cloud storage service trustworthy.

#### *1.2.1.3 Dedicated Cloud network*

Cloud systems (public Clouds specifically) are primarily running based on data centers with dedicated networks, which interconnect with each other using dedicated links [15]. Such a dedicated feature of the Cloud network has provided the Cloud the potential of full bandwidth control ability. The Cloud storage system could benefit from the dedicated Cloud network that the creation and recovery of data can be conducted in a fully controllable and predictable manner. At the meantime, there is still a great potential that data transfer in the Cloud network could be optimized for being conducted more cost-effectively.

#### *1.2.1.4 Big data*

“Big data” is the term for a collection of data sets so large and complex that it becomes difficult to store and process using traditional data storage and processing approaches. In Cloud storage systems, big data is one of the most distinctive features of the Cloud storage system. These data are generated by a large number of Cloud applications, many of which are data and computation intensive and of great importance to these applications. Moreover, the size of the Cloud data is growing even faster. Due to the

huge amount of resources consumed by these data, efficient data management could generate huge value. For managing the massive amounts of Cloud data, the Cloud storage system needs to be powerful enough and able to meet the diverse needs of the data of different usages at different stages.

**1.2.2 The Cloud data life cycle**

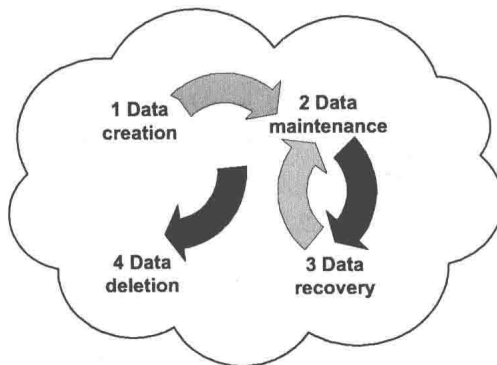
The Cloud data life cycle refers to the period of time starting from the data being created (generated or uploaded) in the Cloud to the data being deleted when the storage space is reclaimed by the Cloud storage system. The life cycle of each piece of Cloud data consists of four stages, which are the data creation stage, the data maintenance stage, the data recovery stage, and the data deletion stage, as depicted in Figure 1.1.

**1.2.2.1 Data creation**

The life cycle of Cloud data starts from the creation of the data in the Cloud storage system. When the original piece of Cloud data (the original replica for short) is created, certain numbers of additional replicas of the Cloud data also need to be created according to the specific reliability requirement of each piece of data and the storage policy [8,9]. All these replicas are transferred and stored on specific storage devices in a distributed fashion.

**1.2.2.2 Data maintenance**

After the data are created and stored, the data maintenance stage commences, which occupies the majority of the Cloud data life. At this stage, Cloud data are processed within applications to achieve different goals. However, for most of the time these data are just stored in storage devices waiting for later use. Certain mechanisms can be conducted to maintain all the replicas so that the service quality is not jeopardized. From the data reliability aspect, the redundancy of Cloud data is maintained at a certain level, so that sufficient data reliability assurance can be offered to meet the storage user's data reliability requirement.



**Figure 1.1 Cloud data life cycle**

### 1.2.2.3 Data recovery

At the data maintenance stage of the Cloud data life cycle, replicas could be lost due to storage failures. In order to either restore the redundancy level of the Cloud data or prevent the data from total loss, data recovery is needed. At this stage, certain mechanisms are conducted to recover the lost replicas. For various purposes, these mechanisms follow different data recovery policies and the duration of the data recovery stage could vary. From the data reliability aspect, the data need to be recovered before the data reliability assurance becomes too low to meet the storage user's requirement.

### 1.2.2.4 Data deletion

When the data are no longer needed, they are deleted. The storage space reclamation mechanism of the Cloud (if any) then recycles the pre-occupied storage space, and the life cycle of the Cloud data ends. Hence this stage of the Cloud data life cycle will not be discussed in this book any further. However, as we will explain later in the book, for determining the proper data reliability assurance that meets the storage user's data reliability requirement, it is preferable that the expected storage duration be given when the data are created.

## 1.3 Key issues of research

The research in this book involves two major aspects: cost-effective data storage and data reliability. On one hand, the storage cost highly depends on the redundancy level of the data. By reducing the redundancy of the Cloud data, the storage cost could be reduced proportionally. Due to the massive amount of big data in the Cloud, the storage cost saved can be huge. On the other hand, reducing redundancy also means that the data reliability may be jeopardized, that is, the data cannot survive until they are deleted (or discarded). In order to provide cost-effective data storage while meeting the data reliability requirement of the Cloud storage users throughout the Cloud data life cycle, our research involves the following key issues.

### 1. Data reliability model

First, we need a model to describe Cloud data reliability and Cloud data reliability-related factors, which is essential for the design of the data reliability assurance approach in the Cloud. The data reliability model should be able to describe the reliability of the Cloud data throughout their life cycles, in which the data are stored with different redundancy levels and stored on different storage devices at different stages respectively.

### 2. Determination of the minimum replication

In order to reduce the storage cost in the Cloud, we need to determine the minimum data redundancy level for meeting the data reliability requirement. As will be further explained in Chapter 3, our research focuses on the data reliability issue in the Cloud with a replication-based data storage scheme. Therefore, in order to store the Cloud data in a cost-effective fashion, at the data creation stage of the Cloud data life cycle, the number of replicas created for the Cloud data need to be minimized. Based on the data reliability model, we need an approach that predicts the data reliability under certain given replication levels so that the



minimum replication that needs to be created can be determined. As a direct consequence, the minimum replication can also act as a benchmark, which can be used for evaluating the cost-effectiveness of various replication-based data storage approaches.

### 3. Cost-effective data reliability assurance

In order to maintain the Cloud data with the minimum replication level, a mechanism that is able to create Cloud data based on the minimum replication calculation approach as well as maintain the created replicas in the Cloud needs to be designed. For effective Cloud data reliability management, this mechanism needs to be able to maintain the big data in the Cloud with a wide variety of data reliability assurance so that all different levels of data reliability requirements can be met. In addition, as a very important aspect, the overhead of such a mechanism also needs to be taken into account.

### 4. Cost-effective data transfer

When replicas of the Cloud data need to be created or are lost, we need to provide effective data transfer process that could maintain the replication level of the data in a cost-effective fashion. In the data creation and recovery stages of the Cloud data life cycle, data transfer activity plays the major role, which transfers the data to the appropriate storage devices. Therefore, optimizing the data transfer in Cloud network could be a good solution for cost-effectiveness. By optimizing data transfer, the cost incurred by data creation or recovery can be reduced.

## 1.4 Book overview

This book systematically investigates the challenging issue of providing cost-effective data storage with data reliability assurance, which includes solid theorems and practical algorithms and finally forms a comprehensive solution to deal with the issue. The book structure is depicted in Figure 1.2.

In Chapter 2, we introduce existing works in literature related to our research. To facilitate our research, literature in three major fields are reviewed. First, from the hardware aspect, to investigate the reliability pattern of storage devices in the Cloud, literature on hardware reliability theories are reviewed. Second, from the software aspect, to investigate data reliability models, and data redundancy maintenance approaches in the Cloud, literature on data reliability modeling, data reliability assurance approaches in distributed data storage systems are reviewed. Third, to investigate data recovery approaches in the Cloud, literature on data recovery and data transfer approaches in distributed systems are reviewed.

In Chapter 3, we present the motivating example of this book and analyze our research problem. We first introduce the motivating example of our research, which is a real-world scientific application for pulsar searching survey of typical data-intensive characteristics. Based on the motivating example, we analyze the research problem and identify details of our research issues.

In Chapter 4, we present our data reliability model for Cloud data storage. Based on the details of our research issues identified in Chapter 3, first we further determine several properties for our data reliability model, and then our novel generic replication-based data reliability model is presented in detail.

In Chapter 5, we present the minimum replication calculation approach. Based on our generic data reliability model presented in Chapter 4, a minimum replication