

Fall 2013

# Efficient integrity verification of replicated data in cloud

Raghul Mukundan

Follow this and additional works at: [http://scholarsmine.mst.edu/masters\\_theses](http://scholarsmine.mst.edu/masters_theses)



Part of the [Computer Sciences Commons](#)

**Department: Computer Science**

---

## Recommended Citation

Mukundan, Raghul, "Efficient integrity verification of replicated data in cloud" (2013). *Masters Theses*. 7201.  
[http://scholarsmine.mst.edu/masters\\_theses/7201](http://scholarsmine.mst.edu/masters_theses/7201)

This Thesis - Open Access is brought to you for free and open access by Scholars' Mine. It has been accepted for inclusion in Masters Theses by an authorized administrator of Scholars' Mine. This work is protected by U. S. Copyright Law. Unauthorized use including reproduction for redistribution requires the permission of the copyright holder. For more information, please contact [scholarsmine@mst.edu](mailto:scholarsmine@mst.edu).



EFFICIENT INTEGRITY VERIFICATION OF REPLICATED DATA IN  
CLOUD

by

RAGHUL MUKUNDAN

A THESIS

Presented to the Faculty of the Graduate School of the  
MISSOURI UNIVERSITY OF SCIENCE AND TECHNOLOGY

In Partial Fulfillment of the Requirements for the Degree

MASTER OF SCIENCE IN COMPUTER SCIENCE

2013

Approved by

Dr. Sanjay Madria, Advisor

Dr. Sriram Chellappan

Dr. Wei Jiang

## PUBLICATION THESIS OPTION

This thesis consists of the following articles that have been submitted for publication as follows:

Paper 1. Pages 10-28 have been published as *Replicated Data Integrity Verification in Cloud*, IEEE Data Engineering Bulletin, December 2012

Paper 2. Pages 29-67 have been submitted as *Efficient Integrity Verification of Replicated Data in Cloud using Homomorphic Encryption*, Distributed and Parallel Database Journal, 2013

## ABSTRACT

The cloud computing is an emerging model in which computing infrastructure resources are provided as a service over the Internet. Data owners can outsource their data by remotely storing them in the cloud and enjoy on-demand high quality services from a shared pool of configurable computing resources. By using these data storage services, the data owners can relieve the burden of local data storage and maintenance. However, since data owners and the cloud servers are not in the same trusted domain, the outsourced data may be at risk as the cloud server may no longer be fully trusted. Therefore, data integrity is of critical importance in such a scenario. Cloud should let the owners or a trusted third party to check for the integrity of their data storage without demanding a local copy of the data. Owners often replicate their data on the cloud servers across multiple data centers to provide a higher level of scalability, availability, and durability. When the data owners ask the Cloud Service Provider (CSP) to replicate data, they are charged a higher storage fee by the CSP. Therefore, the data owners need to be strongly convinced that the CSP is storing data copies agreed on in the service level contract, and data-updates have been correctly executed on all the remotely stored copies. In this thesis, a Dynamic Multi-Replica Provable Data Possession scheme (DMR-PDP) is proposed that prevents the CSP from cheating; for example, by maintaining fewer copies than paid for and/or tampering data. In addition, we also extended the scheme to support a basic file versioning system where only the difference between the original file and the updated file is propagated rather than the propagation of operations for privacy reasons. DMR-PDP also supports efficient dynamic operations like block modification, insertion and deletion on replicas over the cloud servers.

## ACKNOWLEDGMENTS

I owe a debt of gratitude to all those who have helped me with this thesis. First of all, I would like to thank my advisor Dr. Sanjay Kumar Madria, who gave me an opportunity to work on this research. His suggestions and encouragement carried me through difficult times. His valuable feedback contributed greatly to this thesis. Secondly, I would also like to express my gratitude to Dr. Chellappan and Dr. Wei Jiang for serving on my thesis committee and for taking time to review this work.

Im grateful to my loving dad Mukundan Srinivasan for his professional help and inputs. This project is partially funded by the Air Force Research Laboratory (AFRL) in Rome, New York and I express my gratitude to them as well.

## TABLE OF CONTENTS

	Page
PUBLICATION THESIS OPTION .....	ii
ABSTRACT .....	iii
ACKNOWLEDGMENTS .....	iv
LIST OF ILLUSTRATIONS .....	ix
LIST OF TABLES .....	x
 SECTION	
1 INTRODUCTION .....	1
1.1 DATA REPLICATION AND CHALLENGES .....	1
1.2 MOTIVATION .....	2
2 RELATED WORK .....	3
2.1 BASIC APPROACH .....	3
2.2 HOMOMORPHIC VERIFIABLE TAGS .....	3
2.3 PRIVACY-PRESERVING PDP SCHEME .....	4
2.4 MULTIPLE REPLICA PROVABLE DATA POSSESSION SCHEME .....	5
2.5 DYNAMIC PROVABLE DATA POSSESSION .....	5
2.6 ON VERIFYING DYNAMIC MULTIPLE DATA COPIES OVER CLOUD SERVERS .....	6
 PAPER	
I REPLICATED DATA INTEGRITY VERIFICATION IN CLOUD ..	9

ABSTRACT . . . . .	9
1.1 INTRODUCTION . . . . .	10
1.2 Related Work . . . . .	12
1.3 Dynamic Multi-Replica Provable Data Possession (DMR-PDP) Scheme . . . . .	13
1.3.1 Problem Definition and Design Goals . . . . .	14
1.3.2 Preliminaries and Notations . . . . .	15
1.3.3 DMR-PDP Construction . . . . .	16
1.4 Conclusions and Future Work . . . . .	23
BIBLIOGRAPHY . . . . .	24
II EFFICIENT INTEGRITY VERIFICATION OF REPLICATED DATA IN CLOUD USING HOMOMORPHIC ENCRYPTION . . . . .	
ABSTRACT . . . . .	26
2.1 INTRODUCTION . . . . .	27
2.2 Related Work . . . . .	30
2.3 Dynamic Multi-Replica Provable Data Possession (DMR-PDP) Scheme . . . . .	32
2.3.1 Problem Definition and Design Goals . . . . .	32
2.3.2 Preliminaries and Notations . . . . .	33
2.3.3 DMR-PDP Construction . . . . .	35
2.3.4 Using RSA Signatures . . . . .	42
2.4 Security Analysis . . . . .	42
2.5 Multiple Replica File Version Control System (MRFVCS) . . . . .	44
2.5.1 MRFVCS Construction . . . . .	44
2.6 Implementation and Experimental Results . . . . .	50



2.7 Conclusion . . . . .	59
BIBLIOGRAPHY . . . . .	60
SECTION	
3 CONCLUSION . . . . .	62
BIBLIOGRAPHY . . . . .	63
VITA . . . . .	64

## LIST OF ILLUSTRATIONS

Figure	Page
2.1 Rank Based Skip List . . . . .	6
2.2 Merkle Hash Tree . . . . .	7
2.3 MHT Directory . . . . .	8
 PAPER I	
1.1 Cloud Computing Data Storage Model. . . . .	14
1.2 DMR-PDP Scheme . . . . .	17
1.3 Block modification operation in the DMR-PDP scheme . . . . .	22
1.4 Block insertion operation in the DMR-PDP scheme . . . . .	23
 PAPER II	
2.1 Cloud Computing Data Storage Model. . . . .	32
2.2 DMR-PDP Scheme . . . . .	36
2.3 Block modification operation in the DMR-PDP scheme . . . . .	40
2.4 Block insertion operation in the DMR-PDP scheme . . . . .	42
2.5 DMR-PDP scheme using RSA signatures . . . . .	43
2.6 Computation time comparison . . . . .	51
2.7 CSP computation time comparison for the number of replicas on the local cloud servers . . . . .	52
2.8 User computation time comparison for the number of replicas on the local cloud servers . . . . .	53
2.9 Time for file block insert and modify operations on the local cloud servers . . . . .	54
2.10 CSP computation time for number of replicas on Amazon EC2 in- stances . . . . .	55
2.11 Time for file block insert and modify operations in Amazon EC2 micro instance . . . . .	56
2.12 Time for file block insert and modify operations in Amazon EC2 large instance . . . . .	57

2.13 Download and Upload time to EC2 Micro and Large Instances . . .	58
2.14 CSP computation time comparison for FileVersionDeliver algorithm	58

## LIST OF TABLES

Table	Page
PAPER II	
2.1 DMR-PDP Communication cost . . . . .	50

## 1. INTRODUCTION

Cloud computing is an innovative and enticing technology by which highly scalable and technology enabled services can be easily consumed over the Internet. The advantages of cloud computing include on-demand self-service, ubiquitous network access, location independent resource pooling, usage-based pricing, etc. Its great exibility and economic savings are motivating both individuals and enterprises e.g., Amazon, Google, Microsoft, Yahoo, and Salesforce to outsource their data into the cloud. But concerns on data availability and security are preventing data owners and companies from taking advantage of the cloud. When users store data in the cloud, their main concern is whether the cloud can maintain data integrity and whether the data can be recovered when there is data loss or server failure. Cloud Service Providers (CSP), in order to save storage cost, may tend to discard some data or data copies that are not accessed often, or mitigate such data to the second-level storage devices. CSPs may also conceal data loss due to management faults, hardware failures or attacks. Therefore, a critical issue in storing data at untrusted CSPs is periodically verifying whether the storage servers maintain data integrity and store data completely and correctly as stated in the Service Level Agreement (SLA).

### 1.1. DATA REPLICATION AND CHALLENGES

Data Replication is a commonly used technique to increase the data availability in the cloud computing. Cloud replicates the data and stores them strategically on multiple servers located at various geographic locations. Since the replicated data are copies, it is difficult to verify whether the cloud really stores multiple copies of the data. The cloud can easily cheat the owner by storing only one copy of the data. Thus, the owner would like to verify at regular intervals whether the cloud indeed possesses multiple copies of the data as claimed in the SLA. In

general, the cloud has the capability to generate multiple replicas when a data owner challenges the CSP to prove that it possesses multiple copies of the data. Also, it is a valid assumption that the owner of the data may not have a copy of the data stored locally. So, the major task of the owner is not only to verify that the data is intact but also to recover the data if any deletions/corruptions of data are identified. If the owner, during his verification, detects some data loss in any of the replicas in the cloud, he can recover the data from other replicas that are stored intact. Since, the replicas are to be stored at diverse geographic locations, it is assumed to be safe that the data loss will not occur at all the replicas at the same time.

## 1.2. MOTIVATION

Data owners outsource their data to the cloud and may not physically possess a local copy of their data. The data owners should be able to verify that the data with the data stored at the CSP is intact and secure. Cloud should let either users or a user trusted third party to audit the cloud data storage without demanding the local copy of data. These are the factors that need to be considered while designing a scheme that verifies data integrity.

- Communication complexity. The amount of communication between client and server should be minimal.
- Storage cost. Additional storage of client and server required by the scheme should be minimal apart from the original data.
- Data recovery. The scheme should help in data recovery in case of data loss.
- Provable security. The scheme should be secure.

## 2. RELATED WORK

In this section, a review of significant works in the literature on data integrity verification is presented.

### 2.1. BASIC APPROACH

One basic approach to verify data integrity is to hash the entire file before storing it on the cloud. The hash is computed by the data owner and is stored locally. When client wants to verify, he downloads the entire file, computes the hash, and verifies if the values match. But this scheme has the drawbacks of downloading and hashing the entire file every time. This scheme is effective if the file size is less and becomes an overhead for files of huge sizes.

There are two types of schemes that are discussed extensively in the literature, which offer better solution than the basic approach. They are Provable Data Possession (PDP) [1] and Proof of Retrievability (POR) [4]. PDP schemes only checks whether the data stored in the CSP is intact whereas POR schemes help the data owner to recover the data in case of data failure.

### 2.2. HOMOMORPHIC VERIFIABLE TAGS

Ateniese et al [1] proposed this scheme data owner preprocesses the file before storing it on the cloud. They introduce the concept of Homomorphic Verifiable Tags (HVTs). Tags generation is based on RSA signature scheme. Consider a file  $F$  to be a finite ordered collection of  $n$  blocks,  $F = (m_1, m_2, \dots, m_n)$ . Given a message  $m$ , the client computes tag  $T_m$ . These tags will be stored on the cloud together with file  $F$ . These homomorphic verifiable tags act as verification metadata for the file blocks. Because of the homomorphic property, tags computed for multiple file blocks can be combined into a single value. Given two values  $T_{m_i}$  and  $T_{m_j}$ , anyone can combine them into a value  $T_{m_i+m_j}$  corresponding to the sum

of the messages  $m_i + m_j$ . At a later time, the client can verify that the server possesses the file by generating a random challenge against a randomly selected set of file blocks. Using the queried blocks and their corresponding tags, the server generates a proof of possession. The client is thus convinced of data possession, without actually having to retrieve file blocks. The disadvantages of this scheme are that, it does not consider data encryption and data replication. This scheme works only for static files.

### 2.3. PRIVACY-PRESERVING PDP SCHEME

Shah et al.[11] proposed privacy-preserving PDP protocols. Using this scheme, an external Third Party Auditor (TPA) can verify the integrity of files stored by a remote server without knowing any of the file contents. The data owner first encrypts the file, and then sends both the encrypted file along with the encryption key to the remote server. Moreover, the data owner sends the encrypted file along with a key-commitment that fixes a value for the key without revealing the key to the TPA. The primary purpose of this scheme is to ensure that the remote server is correctly possessing clients data along with the encryption key, and to prevent any information leakage to the TPA which is responsible for the auditing task. Since key should be kept secret from both the auditor and the cloud, the data owner places a key commitment  $gk$  instead of storing  $k$  on the server. The auditor generates  $n$  random numbers and generates MAC values for random blocks and stores it. Periodically the auditor issues a challenge to the server by generating a random block number for which the server generates MAC for that block and sends it back to the auditor for verification. This scheme has a lot of disadvantages. The number of times a particular data item can be verified is limited and must be fixed beforehand. TPA has to regenerate a new list of hash values to achieve unbounded number of audits. Lack of support for stateless verification, i.e., the TPA has to update its state between audits to prevent using the same random number or the same MAC twice.



## 2.4. MULTIPLE REPLICA PROVABLE DATA POSSESSION SCHEME

Curtmola [6] proposed Multiple-Replica PDP (MR-PDP) scheme where a data owner can verify that several copies of a file are stored by a storage service provider. The MR-PDP scheme is an extension to the PDP models proposed by Ateniese. Curtmola proposed creating distinct replicas/copies of the data file by first encrypting the file then masking the encrypted version with some randomness generated from a Pseudo-Random Function (PRF). Different PRF keys are used to generate multiple data copies. RSA signatures are used for tag generation. RSA signatures have holomorphic property where multiple data blocks can be verified at the same time. The main advantage of this scheme is that it proves the integrity of multiple replicas. The disadvantages are the data is not encrypted, the size of RSA signatures is huge, deals with static data, and the authorized users have to know which copy has been specifically retrieved from the CSP to properly unmask it before decryption.

## 2.5. DYNAMIC PROVABLE DATA POSSESSION

This scheme [7] tries to solve the certain limitations of MRPDP scheme by using a rank based authenticated skip list. This data structure is used to verify the integrity of the data. According to this scheme, a file  $F$  is split into  $n$  blocks  $m_1, m_2, \dots, m_n$ . The tag  $T(m_i)$  of block  $m_i$  is stored at the  $i$ -th bottom-level node of the skip list. Block  $m_i$  will be stored elsewhere by the cloud. Each node  $v$  of the skip list stores the number of nodes at the bottom level that can be reached from  $v$ . This value is called the rank of  $v$  and denote it with  $r(v)$ . The figure below shows a skip list with ranks of nodes. An insertion, deletion, or modification of a file block affects only the nodes of the skip list along a search path. Ranks of the affected nodes are computed bottom-up in constant time.

The top leftmost node of a skip list will be referred to as the start node. For a node  $v$  the indices of the leftmost and rightmost nodes at the bottom level

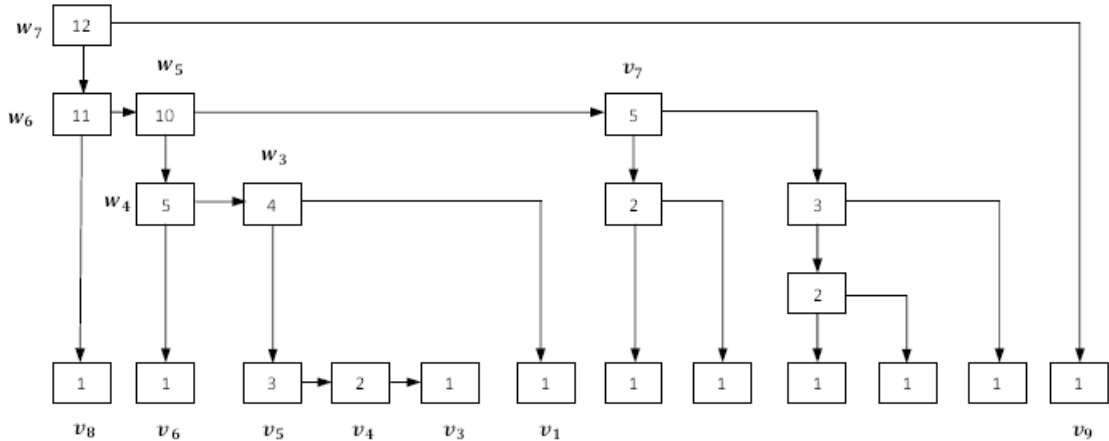


Figure 2.1. Rank Based Skip List

reachable from  $v$ , is denoted by  $\text{low}(v)$  and  $\text{high}(v)$  respectively. At each level based on certain conditions each node will also store the collective hash value of the nodes from left to right in the bottom level in this way  $h(x_1, \dots, x_k) = h(h(x_1) \parallel \dots \parallel h(x_k))$ . Where  $\parallel$  denote a concatenation operator. These hash entries act as metadata. Whenever a file block is updated or modified, the hash value of the nodes in this skip list from which the updated or modified node is reachable is recomputed and the path where the nodes are updated is sent to the user for verification. The integrity verification scheme is similar to the earlier schemes but this scheme supports data dynamics and verifies data block insert, update and modify operations.

## 2.6. ON VERIFYING DYNAMIC MULTIPLE DATA COPIES OVER CLOUD SERVERS

This scheme [8] supports multiple replicas and data encryption. The data blocks in each copy are appended to file blocks before encryption. AES encryption scheme is used for data encryption and BLS signatures are used for tag generation. Dynamic data operations are supported with the help of Merkle Hash Trees (MHT).

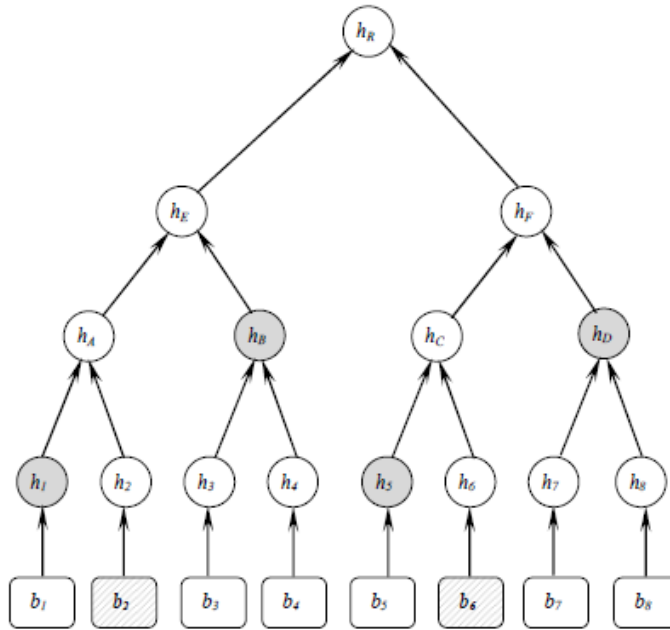


Figure 2.2. Merkle Hash Tree

A MHT is a binary tree structure used to efficiently verify the integrity of the data. The MHT is a tree of hashes where leaves of the tree are the hashes of the data blocks. In the above figure  $h_A = h(h_1 \parallel h_2)$ ,  $h_B = h(h_3 \parallel h_4)$ , and so on. Finally,  $h_R = h(h_E \parallel h_F)$  is the hash of the root node that is used to authenticate the integrity of all data blocks. The data blocks  $\{b_1, b_2, \dots, b_8\}$  are stored on a remote server, and only the authentic value  $h_R$  is stored locally on the verifier side. For example, if the verifier requests to check the integrity of the blocks  $b_2$  and  $b_6$ , the server will send these two blocks along with the authentication paths  $A_2 = \{h_1, h_B\}$  and  $A_6 = \{h_5, h_D\}$  that are used to reconstruct the root of the MHT.  $A_j$  the authentication path of  $b_j$  is a set of node siblings (grey-shaded circles) on the path from  $h_j$  to the root. In the dynamic behavior of outsourced data, both the values and the positions of the data block needs to be authenticated. This will give an assurance that a specific value is stored at a specific leaf node. For example, if a data owner requires inserting a new block after position  $j$ , the verifier needs to make sure that the server has inserted the new block in the requested position. To validate the positions of the blocks, the leaf nodes of the MHT are

treated in a specific sequence, e.g., left-to-right sequence. So, the hash of any internal node is  $h(\text{left child} \parallel \text{right child})$ , e.g.,  $h_A = h(h_1 \parallel h_2) \neq h(h_2 \parallel h_1)$ . Besides, the authentication path  $A_j$  is viewed as an ordered set, and thus any leaf node is uniquely specified by following the used sequence of constructing the root of the MHT. A MHT directory is constructed for all the replicas.

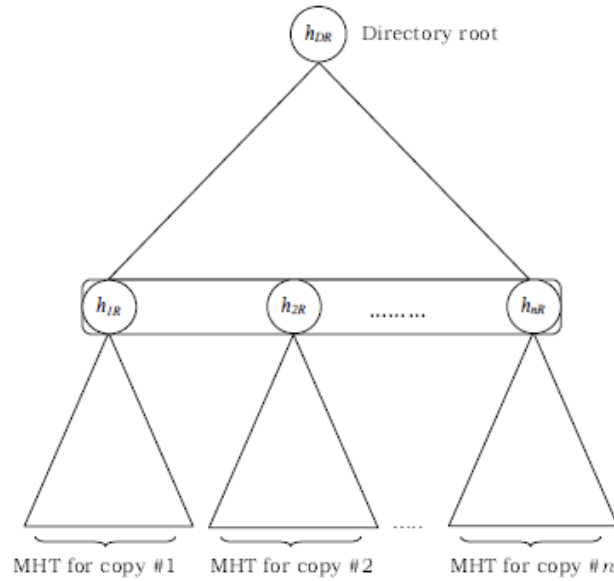


Figure 2.3. MHT Directory

The disadvantage of this scheme is that the number of nodes in MHT depends on number of blocks in the file and also on the number of copies. If a file has huge number of blocks, then number of nodes in MHT will be huge. This scheme incurs computation and communication overhead to the data owner to generate and send such a huge tree structure. For every verification, the authentication paths and multiple hash values are also sent to the user which will create a communication overhead to data owner.

**PAPER I****REPLICATED DATA INTEGRITY VERIFICATION IN CLOUD**

Raghul Mukundan<sup>1</sup>, Sanjay Madria<sup>1</sup>, Mark Linderman<sup>2</sup>

<sup>1</sup>*Department of Computer Science, Missouri University of Science & Technology,  
Rolla, MO 65409*

<sup>2</sup>*Air Force Research Lab, Rome, NY*

**ABSTRACT**

Cloud computing is an emerging model in which computing infrastructure resources are provided as a service over the Internet. Data owners can outsource their data by remotely storing them in the cloud and enjoy on-demand high quality applications and services from a shared pool of configurable computing resources. However, since data owners and cloud servers are not in the same trusted domain, the outsourced data may be at risk as the cloud server may no longer be fully trusted. Therefore, data integrity is of critical importance in such a scenario. Cloud should let either the owners or a trusted third party to audit their data storage without demanding a local copy of the data from owners. Replicating data on cloud servers across multiple data centers provides a higher level of scalability, availability, and durability. When the data owners ask the Cloud Service Provider (CSP) to replicate data at different servers, they are charged a higher fee by the CSP. Therefore, the data owners need to be strongly convinced that the CSP is storing all the data copies that are agreed upon in the service level contract, and the data-update requests issued by the customers have been correctly executed on all the remotely stored copies. To deal with such problems, previous multi copy verification schemes either focused on static files or incurred huge update costs in a dynamic file scenario. In this paper, we propose some ideas under a Dynamic

Multi-Replica Provable Data Possession scheme (DMR-PDP) that prevents the CSP from cheating; for example, by maintaining fewer copies than paid for. DMR-PDP also supports efficient dynamic operations like block modification, insertion and deletion on data replicas over cloud servers.

## 1.1. INTRODUCTION

When users store data in the cloud, their main concern is whether the cloud can maintain the data integrity and data can be recovered when there is a data loss or server failure. Cloud service providers (CSP), in order to save storage cost, may tend to discard some data or data copies that is not accessed often, or mitigate data to second-level storage devices. CSPs may also conceal data loss due to management faults, hardware failures or attacks. Therefore, a critical issue in storing data at untrusted CSPs is periodically verifying whether the storage servers maintain data integrity; store data completely and correctly as stated in the service level agreement (SLA).

Data replication is a commonly used technique to increase the data availability in cloud computing. Cloud replicates the data and stores them strategically on multiple servers located at various geographic locations. Since the replicated copies look exactly similar, it is difficult to verify whether the cloud really stores multiple copies of the data. Cloud can easily cheat the owner by storing only one copy of the data. Thus, the owner would like to verify at regular intervals whether the cloud indeed possesses multiple copies of the data as claimed in the SLA. In general, cloud has the capability to generate multiple replicas when a data owner challenges the CSP to prove that it possesses multiple copies of the data. Also, it is a valid assumption that the owner of the data may not have a copy of the data stored locally. So, the major task of the owner is not only to verify that the data is intact but also to recover the data if any deletions/corruptions of data are identified. If the data owner during his verification using DMR-PDP scheme detects some data loss in any of the replicas in the cloud, he can recover the data

from other replicas that are stored intact. Since, the replicas are to be stored at diverse geographic locations, it is safe to assume that a data loss will not occur at all the replicas at the same time.

Provable data possession (PDP) [2] is a technique to audit and validate the integrity of data stored on remote servers. In a typical PDP model, the data owner generates metadata/tag for a data file to be used later for integrity verification. To ensure security, the data owner encrypts the file and generates tags on the encrypted file. The data owner sends the encrypted file and the tags to the cloud, and deletes the local copy of the file. When the data owner wishes to verify the data integrity, he generates a challenge vector and sends it to the cloud. The cloud replies by computing a response on the data and sends it to the verifier/data owner to prove that multiple copies of the data file are stored in the cloud. Different variations of PDP schemes such as [2], [4], [6], [7], [9], [10], [11], [12], [15] were proposed under different cryptographic assumptions. But most of these schemes deal only with static data files and are valid only for verifying a single copy. A few other schemes such as [3], [5], [8], [13], [14] provide dynamic scalability of a single copy of a data file for various applications which mean that the remotely stored data can not only be accessed by the authorized users, but also be updated and scaled by the data owner.

In this paper, we propose a scheme that allows the data owner to securely ensure that the CSP stores multiple replicas. A simple way to make the replicas look unique and differentiable is using probabilistic encryption schemes. Probabilistic encryption creates different cipher texts each time the same message is encrypted using the same key. Thus, our scheme uses homomorphic probabilistic encryption to create distinct replicas/copies of the data file and BLS signatures [17] to create constant amount of metadata for any number of replicas. Probabilistic encryption encrypts all the replicas with the same key. Therefore, in our scheme the data owner will have to share just one decryption key with the authorized users and need not worry about CSP granting access to any of the replicas to

the authorized users. The homomorphic property of the encryption scheme helps in efficient file updates. The data owner has to encrypt the difference between the updated file and the old file and send it to the cloud, which updates all the replicas by performing homomorphic addition on the file copies. Any authenticated data structure, e.g, Merkle Hash Trees and Skiplist can be used with our scheme to ensure that the cloud uses the right file blocks for data integrity verification. However, the ways to efficiently manage authenticated data structures in the cloud is not within the scope of our paper.

Organization: The rest of the paper is organized as follows. Overview of the related work is provided in Section 2 followed by the problem definition in Section 3, and a detailed description of our scheme in Section 4. Future work is discussed in Section 5.

## 1.2. RELATED WORK

Ateniese et al. [2] were the first to define the Provable Data Possession (PDP) model for ensuring the possession of files on untrusted storages. They made use of RSA-based homomorphic tags for auditing outsourced data. However, dynamic data storage and multiple replica system are not considered in this scheme. In their subsequent work [12], they proposed a dynamic version which supports very basic block operations with limited functionality, and does not support block insertions. In [13], Wang et al. considered dynamic data storage in a distributed scenario, and proposed a challenge-response protocol which can determine data correctness as well as locate possible errors. Similar to [12], only partial support for dynamic data operation is considered. Erway et al. [14] extended the PDP model in [2] to support provable updates to stored data files using rank-based authenticated skip lists. However, efficiency of their scheme remains unclear and these schemes hold good only for verifying a single copy.



Curtmola et.al [1] proposed Multiple-Replica PDP (MR-PDP) scheme wherein the data owner can verify that several copies of a file are stored by a storage service provider. In their scheme, distinct replicas are created by first encrypting the data and then masking it with randomness generated from a Pseudo-Random Function (PRF). The randomized data is then stored across multiple servers. The scheme uses RSA signatures for creation of tags. But, their scheme did not address how the authorized users of the data can access the file copies from the cloud servers noting that the internal operations of the CSP are opaque and do not support dynamic data operations. Ayad F. Barsoum et al. [16] proposed creation of distinct copies by appending replica number to the file blocks and encrypting it using an encryption scheme that has strong diffusion property, e.g., AES. Their scheme supports dynamic data operations but during file updates, the copies in all the servers should be encrypted again and updated on the cloud. This scheme suits perfectly for static multiple replicas but proves costly in a dynamic scenario. BLS signatures are used for creating tags and authenticated data structures like Merkle Hash Trees are used to ensure right file blocks are used during verification. Authorized users of the data should know random numbers in [1] and replica number in [16] to generate the original file.

### **1.3. DYNAMIC MULTI-REPLICA PROVABLE DATA POSSESSION (DMR-PDP) SCHEME**

The cloud computing model considered in this work consists of three main components as illustrated in Figure 1: (i) a data owner that can be an individual or an organization originally possessing sensitive data to be stored in the cloud; (ii) a CSP who manages cloud servers and provides paid storage space on its infrastructure to store the owner's files and (iii) authorized users - a set of owner's clients who have the right to access the remote data and share some keys with the data owner.



Figure 1.1. Cloud Computing Data Storage Model.

**1.3.1. Problem Definition and Design Goals.** More recently, many data owners relieve their burden of local data storage and maintenance by outsourcing their data to a CSP. CSP undertakes the data replication task in order to increase the data availability, durability and reliability but the customers have to pay for using the CSPs storage infrastructure. On the other hand, cloud customers should be convinced that the (1) CSP actually possesses all the data copies as agreed upon, (2) integrity of these data copies are maintained, and (3) the customers are receiving the service that they are paying for. Therefore, in this paper, we address the problem of securely and efficiently creating multiple replicas of the data file of the owner to be stored over untrusted CSP and then auditing all these copies to verify their completeness and correctness. Our design goals are summarized below:

1. Dynamic Multi-Replica Provable Data Possession (DMR-PDP) protocols should efficiently and securely provide the owner with strong evidence that the CSP is in possession of all the data copies as agreed upon and that these copies are intact.
2. Allowing the users authorized by the data owner to seamlessly access a file copy from the CSP.

3. Using only a single set of metadata/tags for all the file replicas for verification purposes.
4. Allowing dynamic data operation support to enable the data owner to perform block-level operations on the data files while maintaining the same level of data correctness assurance.
5. Enabling both probabilistic and deterministic verification guarantees.

**1.3.2. Preliminaries and Notations.** In this section, we provide details of the Bilinear mapping and Paillier Encryption schemes used in our present work.

1. Assume that  $F$ , a data file to be outsourced, is composed of a sequence of  $m$  blocks, i.e.,  $F = \{b_1, b_2, \dots, b_m\}$ .
2.  $F_i = \{b_{i1}, b_{i2}, \dots, b_{im}\}$  represents the file copy  $i$ .
3. Bilinear Map/Pairing: Let  $G_1$ ,  $G_2$ , and  $G_T$  be cyclic groups of prime order  $a$ . Let  $u$  and  $v$  be generators of  $G_1$  and  $G_2$ , respectively. A bilinear pairing is a map  $e : G_1 \times G_2 \rightarrow G_T$  with the following properties:
  - Bilinear:  $e(u_1 u_2, v_1) = e(u_1, v_1) \cdot e(u_2, v_1)$ ,  $e(u_1, v_1 v_2) = e(u_1, v_1) \cdot e(u_1, v_2) \forall u_1, u_2 \in G_1$  and  $v_1, v_2 \in G_2$
  - Non-degenerate:  $e(u, v) \neq 1$
  - There exists an efficient algorithm for computing  $e$
  - $e(u_1^x, v_1^y) = e(u_1, v_1)^{xy} \forall u_1 \in G_1; v_1 \in G_2$ , and  $x, y \in \mathbb{Z}_a$
4.  $H(\cdot)$  is a map-to-point hash function:  $\{0, 1\}^* \rightarrow G_1$ .
5. Homomorphic Encryption: A homomorphic encryption scheme has the following properties.
  - $E(m_1 + m_2) = E(m_1) +_h E(m_2)$  where  $+_h$  is a homomorphic addition operation.
  - $E(k * m) = E(m)^k$ .

where  $E(.)$  represents a homomorphic encryption scheme and  $m, m_1, m_2$  are messages that are encrypted and  $k$  is some random number.

6. **Paillier Encryption:** Paillier cryptosystem is a homomorphic probabilistic encryption scheme. The steps are as follows.

- Compute  $N = p * q$  and  $\lambda = \text{LCM}(p-1, q-1)$ , where  $p, q$  are two prime numbers.
- Select a random number  $g$  such that its order is a multiple of  $N$  and  $g \in \mathbb{Z}_N^{2*}$ .
- Public key is  $(N, g)$  and secret key is  $\lambda$ , where  $N = p*q$ .
- Cipher text for a message  $m$  is computed as  $c = g^m r^N \text{ mod } N^2$  where  $r$  is a random number and  $r \in \mathbb{Z}_N^*$ ,  $c \in \mathbb{Z}_N^{2*}$  and  $m \in \mathbb{Z}_N$ .
- Plain text is obtained by  $m = L(c^\lambda \text{ mod } N^2) * (L(g^\lambda \text{ mod } N^2))^{-1} \text{ mod } N$ .

7. Properties of public key  $g$  in Paillier Scheme

- $g \in \mathbb{Z}_N^{2*}$ .
- If  $g = (1 + N) \text{ mod } N^2$ , it has few interesting properties
  - (a) Order of the value  $(1 + N)$  is  $N$ .
  - (b)  $(1 + N)^m \equiv (1 + mN) \text{ mod } N^2$ .  $(1 + mN)$  can be used directly instead of calculating  $(1 + N)^m$ . This avoids the costly exponential operation during data encryption.

**1.3.3. DMR-PDP Construction.** In our approach, the data owner creates multiple encrypted replicas and uploads them on to the cloud. The CSP stores them on one or multiple servers located at various geographic locations. The data owner shares the decryption key with a set of authorized users. In order to access the data, an authorized user sends a data request to the CSP and receives a data copy in an encrypted form that can be decrypted using a secret key shared with

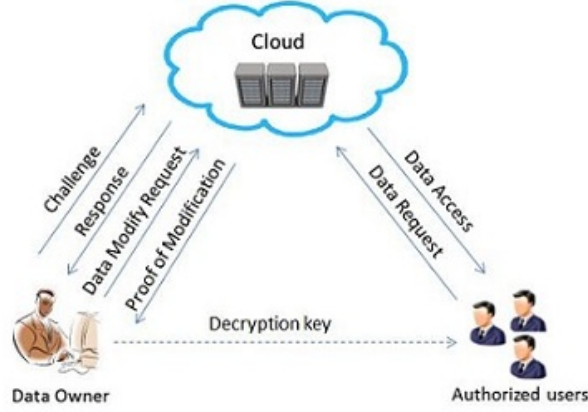


Figure 1.2. DMR-PDP Scheme

the owner. The proposed scheme consists of seven algorithms: KeyGen, ReplicaGen, TagGen, Prove, Verify, PrepareUpdate and ExecUpdate. The overview of the communication involved in our scheme is shown in Figure 2.

1.  $(pk, sk) \leftarrow \text{KeyGen}()$ . This algorithm is run by the data owner to generate a public key  $pk$  and a private key  $sk$ . The data owner generates three sets of keys.
  - (a) Keys for data tags : This key is used for generating tags for the data. The data owner selects a bilinear map  $e$  and selects a private key  $l \in \mathbb{Z}_a$ . Public key is calculated as  $y = v^l \in G_2$ .
  - (b) Keys for data : This key is used for encrypting the data and thereby creating multiple data copies. The data owner selects paillier public keys  $(N, g)$  with  $g = (1 + N) \bmod N^2$  and secret key  $\lambda$ .
  - (c) PRF key : The data owner generates a PRF key  $Key_{PRF}$  which generates  $s$  numbers. These  $s$  numbers are used in creating  $s$  copies of the data. Each number is used in creating one data copy. Let  $\{k_1, k_2, \dots, k_s\} \in \mathbb{Z}_N^*$  be the numbers generated by the PRF key.  $Key_{PRF}$  is maintained confidentially by the data owner and hence the  $s$  numbers used in creating multiple copies are not known to the cloud.
2.  $\{F_i\}_{1 \leq i \leq s} \leftarrow \text{ReplicaGen}(s, F)$ . This algorithm is run by the data owner. It takes the number of replicas  $s$  and the file  $F$  as input and generates  $s$  unique

differentiable copies  $\{F_i\}_{1 \leq i \leq s}$ . This algorithm is run only once. Unique copies of each file block of file  $F$  is created by encrypting it using a probabilistic encryption scheme, e.g., Paillier encryption scheme. Through probabilistic encryption, encrypting a file block  $s$  times yields  $s$  distinct ciphertexts. For a file  $F = \{b_1, b_2, \dots, b_m\}$  multiple data copies are generated using Paillier encryption scheme as  $F_i = \{(1+N)^{b_1}(k_i r_{i1})^N, (1+N)^{b_2}(k_i r_{i2})^N, \dots, (1+N)^{b_m}(k_i r_{im})^N\}_{1 \leq i \leq m}$ . Using Paillier's properties the above result can be rewritten as  $F_i = \{(1+b_1N)(k_i r_{i1})^N, (1+b_2N)(k_i r_{i2})^N, \dots, (1+b_mN)(k_i r_{im})^N\}_{1 \leq i \leq m}$ , where  $i$  represents the file copy number,  $k_i$  represents the numbers generated from PRF key  $Key_{PRF}$  and  $r_{ij}$  represents any random number used in Paillier encryption scheme.  $k_i$  is multiplied by a random number  $r_{ij}$  and the product is used for encryption. The presence of  $k_i$  in a file block identifies which file copy the file block belongs to. All these file copies yield the original file when decrypted. This allows the users authorized by the data owner to seamlessly access the file copy received from the CSP.

3.  $\phi \leftarrow \text{TagGen}(\text{sk}, F)$ . This algorithm is run by the data owner. It takes the private key  $\text{sk}$  and the file  $F$  as input and outputs the tags  $\phi$ . We use BLS signature scheme to create tags on the data. BLS signatures are short and homomorphic in nature and allow concurrent data verification, which means multiple data blocks can be verified at the same time. In our scheme, tags are generated on each file block  $b_i$  as  $\phi_i = (H(F) \cdot u^{b_i N})^l \in G_1$  where  $u \in G_1$  and  $H(\cdot) \in G_1$  represents hash value which uniquely represents the file  $F$ . The data owner sends the tag set  $\phi = \{\phi_i\}_{1 \leq i \leq m}$  to the cloud.
4.  $P \leftarrow \text{Prove}(F, \phi, \text{challenge})$ . This algorithm is run by the CSP. It takes the file replicas of file  $F$ , the tags  $\phi$  and *challenge* vector sent by the data owner as input and returns a proof  $P$  which guarantees that the CSP is actually storing  $s$  copies of the file  $F$  and all these copies are intact. The data owner uses the proof  $P$  to verify the data integrity. There are two phases in this algorithm:

(a) **Challenge:** In this phase the data owner challenges the cloud to verify the integrity of all outsourced copies. There are two types of verification schemes:

- i. Deterministic - here all the file blocks from all the copies are used for verification.
- ii. Probabilistic - only a few blocks from all the copies are used for verification. A Pseudo Random Function key (PRF) is used to generate random indices ranging between 1 and  $m$ . The file blocks from these indices are used for verification. In each verification a percentage of the file is verified and it accounts for the verification of the entire file.

At each challenge, the data owner chooses the type of verification scheme he wishes to use. If the owner chooses the deterministic verification scheme, he generates one PRF key,  $Key_1$ . If he chooses the probabilistic scheme he generates two PRF keys,  $Key_1$  and  $Key_2$ . PRF keyed with  $Key_1$  generates  $c$  ( $1 \leq c \leq m$ ) random file indices which indicates the file blocks that CSP should use for verification. PRF keyed with  $Key_2$  generates  $s$  random values and the CSP should use each of these random numbers for each file copy while computing the response. The data owner sends the generated keys to the CSP.

(b) **Response:** This phase is executed by the CSP when a challenge for data integrity verification is received from the data owner. Here, we show the proof for probabilistic verification scheme (the deterministic verification scheme also follows the same procedure). The CSP receives two PRF keys,  $Key_1$  and  $Key_2$  from the data owner. Using  $Key_1$ , CSP generates a set  $\{C\}$  with  $c$  ( $1 \leq c \leq m$ ) random file indices ( $\{C\} \in \{1, 2, \dots, m\}$ ), which indicate the file blocks that CSP should use for verification. Using  $Key_2$ , CSP generates 's' random values  $T = \{t_1,$

$t_2, \dots, t_s$ . The cloud performs two operations. One on the tags and the other on the file blocks.

- i. Operation on the tags: Cloud multiplies the file tags corresponding to the file indices generated by PRF key  $\text{Key}_1$ .

$$\begin{aligned}
 \sigma &= \prod_{j \in C} (H(F) \cdot u^{b_j N})^l \\
 &= \prod_{j \in C} H(F)^l \cdot \prod_{j \in C} u^{b_j N l} \\
 &= H(F)^{cl} \cdot u^{Nl \sum_{j \in C} (b_j)}
 \end{aligned}$$

- ii. Operation on the file blocks: The cloud first takes each file copy and multiplies all the file blocks corresponding to the file indices generated by PRF key  $\text{Key}_1$ . The product of each copy is raised to the power the random number generated for that copy by the PRF key  $\text{Key}_2$ . The result of the above operation for each file copy  $i$  is given by  $(\prod_{j \in C} (1+N)^{b_j} (k_i r_{ij})^N)^{t_i} \bmod N^2$ . The CSP then multiplies the result of each copy to get the result

$$\begin{aligned}
 \mu &= \prod_{i=1}^s \left( \prod_{j \in C} (1+N)^{b_j} (k_i r_{ij})^N \right)^{t_i} \\
 &= \prod_{i=1}^s \left( \prod_{j \in C} (1+N)^{b_j t_i} \prod_{j \in C} (k_i r_{ij})^{N t_i} \right) \\
 &= \prod_{i=1}^s \left( (1+N)^{t_i \sum_{j \in C} b_j} \prod_{j \in C} (k_i r_{ij})^{N t_i} \right) \\
 &= \left( \prod_{i=1}^s (1+N)^{t_i \sum_{j \in C} b_j} \right) \left( \prod_{i=1}^s \left( (k_i)^{c t_i N} \prod_{j \in C} (r_{ij})^{N t_i} \right) \right) \\
 &= \left( (1+N)^{\sum_{i=1}^s t_i \sum_{j \in C} b_j} \right) \left( \prod_{i=1}^s (k_i)^{c t_i N} \right) \left( \prod_{i=1}^s \prod_{j \in C} (r_{ij})^{N t_i} \right)
 \end{aligned}$$



Using properties of Paillier scheme, the above equation can be rewritten as

$$\mu = (1 + N \sum_{i=1}^s (t_i) \sum_{j \in C} (b_j)) (\prod_{i=1}^s (k_i)^{Nct_i}) (\prod_{i=1}^s (\prod_{j \in C} (r_{ij})^{t_i N}))$$

The CSP sends  $\sigma$  and  $\mu \bmod N^2$  values to the data owner.

5.  $\{1, 0\} \leftarrow \text{Verify}(\text{pk}, \text{P})$ . This algorithm is run by the data owner. It takes as input the public key  $\text{pk}$  and the proof  $\text{P}$  returned from the CSP, and outputs 1 if the integrity of all file copies is correctly verified or 0 otherwise. After receiving  $\sigma$  and  $\mu$  values from the CSP, the data owner does the following:
  - (a) calculates  $v = (\prod_{i=1}^s (k_i)^{t_i c N})$  and  $d = \text{Decrypt}(\mu) / (\sum_{i=1}^s k_i)$ . This can be calculated from the values generated from  $\text{Key}_{PRF}$  and  $c$ .
  - (b) checks if  $\mu \bmod v \equiv 0$ . This ensures that the cloud has used all the file copies while computing the response.
  - (c) checks if  $(\text{H}(\text{F})^c u^{dN})^l = \sigma$ . This ensures that the CSP has used all the file blocks while computing the response. If options b and c are satisfied, it indicates that the data stored by the owner in the cloud is intact and the cloud has stored multiple copies of the data as agreed in the service level agreement.
6. *Update*  $\leftarrow \text{PrepareUpdate}()$ . This algorithm is run by the data owner to perform any operation on the outsourced file copies stored by the remote CSP. The output of this algorithm is an *Update* request. The data owner sends the *Update* request to the cloud and will be of the form  $\langle \text{Id}_F, \text{BlockOp}, j, b_i', \phi' \rangle$ , where  $\text{Id}_F$  is the file identifier,  $\text{BlockOp}$  corresponds to block operation,  $j$  denotes the index of the file block,  $b_i'$  represents the updated file blocks and  $\phi'$  is the updated tag.  $\text{BlockOp}$  can be data modification, insertion or delete operation.

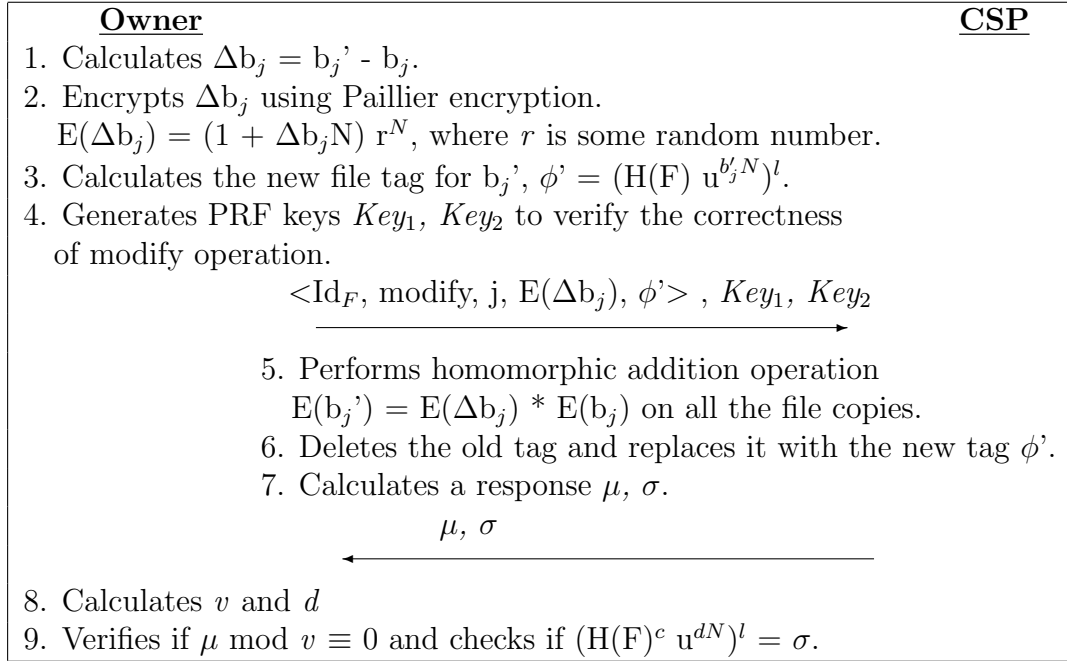


Figure 1.3. Block modification operation in the DMR-PDP scheme

7.  $(F', \phi') \leftarrow \text{ExecUpdate}(F, \phi, \text{Update})$ . This algorithm is run by the CSP where the input parameters are the file copies  $F$ , the tags  $\phi$ , and  $\text{Update}$  request (sent from the owner). It outputs an updated version of all the file copies  $F'$  along with updated signatures  $\phi'$ . After any block operation, the data owner runs the challenge protocol to ensure that the cloud has executed the operations correctly. The operation in  $\text{Update}$  request can be modifying a file block, inserting a new file block or deleting a file block.

(a) **Modification:** Data modification is one of the most frequently used dynamic operations. The data modification operation in DMR-PDP scheme is shown in Figure 3.

(b) **Insertion:** In the block insertion operation, the owner inserts a new block after position  $j$  in a file. If the file  $F$  had  $m$  blocks initially, the file will have  $m+1$  blocks after the insert operation. The file block insertion operation is shown in Figure 4.

(c) **Deletion:** Block deletion operation is the opposite of the insertion operation. When one block is deleted, indices of all subsequent blocks

are moved one step forward. To delete a specific data block at position  $j$  from all copies, the owner sends a delete request  $\langle \text{Id}_F, \text{delete}, j, \text{null}, \text{null} \rangle$  to the cloud. Upon receiving the request, the cloud deletes the tag and the file block at index  $j$  in all the file copies.

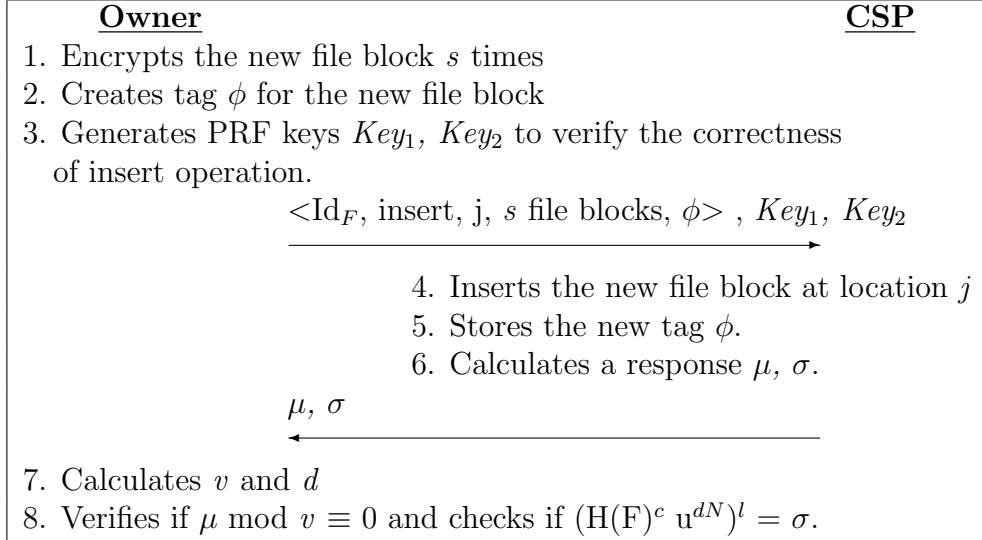


Figure 1.4. Block insertion operation in the DMR-PDP scheme

#### 1.4. CONCLUSIONS AND FUTURE WORK

In this paper, we have discussed work related to the replicated data integrity preservation in a cloud environment and presented a Dynamic Multi-Replica Provable Data Possession scheme (DMR-PDP) to periodically verify the correctness and completeness of multiple data copies stored in the cloud. Our scheme also supports dynamic data update operations. All the data copies can be decrypted using a single decryption key, thus providing a seamless access to the data's authorized users. This scheme can be extended for multiple versions where only deltas can be stored in the cloud and owner can save on storage cost. Currently, we are implementing the proposed scheme for evaluating it in a real cloud platform using different performance metrics and comparing it with some of the existing methods. We also plan to extend this scheme for secure multi-version data where only one original and multiple deltas can be stored in the cloud.

**BIBLIOGRAPHY**

- [1] R. Curtmola, O. Khan, R. Burns, and G. Ateniese, “MR-PDP: Multiple-Replica Provable Data Possession,” in 28th IEEE ICDCS, 2008, pp. 411-420.
- [2] G. Ateniese, R. Burns, R. Curtmola, J. Herring, L. Kissner, Z. Peterson, and D. Song, “Provable data possession at untrusted stores,” in CCS ’07: Proceedings of the 14th ACM Conference on Computer and Communications Security, New York, NY, USA, 2007, pp. 598-609.
- [3] G. Ateniese, R. D. Pietro, L. V. Mancin, and G. Tsudik, “Scalable and efficient provable data possession,” in SecureComm 08: Proceedings of the 4th International Conference on Security and Privacy in Communication Networks, New York, NY, USA, 2008, pp. 1-10.
- [4] Y. Deswarte, J.-J. Quisquater, and A. Sadane, “Remote integrity checking,” in 6th Working Conference on Integrity and Internal Control in Information Systems (IICIS), S. J. L. Strous, Ed., 2003, pp. 1-11.
- [5] C. Erway, A. Kupcu, C. Papamanthou, and R. Tamassia, “Dynamic provable data possession,” in CCS 09: Proceedings of the 16th ACM Conference on Computer and Communications Security, New York, NY, USA, 2009, pp. 213-222.
- [6] D. L. G. Filho and P. S. L. M. Barreto, “Demonstrating data possession and uncheatable data transfer,” Cryptology ePrint Archive, Report 2006/150, 2006.
- [7] P. Golle, S. Jarecki, and I. Mironov, “Cryptographic primitives enforcing communication and storage complexity,” in FC’02: Proceedings of the 6th International Conference on Financial Cryptography, Berlin, Heidelberg, 2003, pp. 120-135.
- [8] Z. Hao, S. Zhong, and N. Yu, “A privacy-preserving remote data integrity checking protocol with data dynamics and public verifiability,” IEEE Transactions on Knowledge and Data Engineering, vol. 99, no. PrePrints, 2011.
- [9] E. Mykletun, M. Narasimha, and G. Tsudik, “Authentication and integrity in outsourced databases,” Trans. Storage, vol. 2, no. 2, 2006.
- [10] F. Sebe, J. Domingo-Ferrer, A. Martinez-Balleste, Y. Deswarte, and J.-J. Quisquater, “Efficient remote data possession checking in critical information infrastructures,” IEEE Trans. on Knowl. and Data Eng., vol. 20, no. 8, 2008.
- [11] M. A. Shah, M. Baker, J. C. Mogul, and R. Swaminathan, “Auditing to keep online storage services honest,” in HOTOS’07: Proceedings of the 11th USENIX workshop on Hot topics in operating systems, Berkeley, CA, USA, 2007, pp. 1-6.

- [12] M. A. Shah, R. Swaminathan, and M. Baker, “Privacy-preserving audit and extraction of digital contents,” Cryptology ePrint Archive, Report 2008/186, 2008.
- [13] C. Wang, Q. Wang, K. Ren, and W. Lou, “Ensuring data storage security in cloud computing,” Cryptology ePrint Archive, Report 2009/081, 2009, <http://eprint.iacr.org> (Retrieved: 08/06/2013).
- [14] Q. Wang, C. Wang, J. Li, K. Ren, and W. Lou, “Enabling public verifiability and data dynamics for storage security in cloud computing,” in ESORICS09: Proceedings of the 14th European Conference on Research in Computer Security, Berlin, Heidelberg, 2009, pp. 355-370.
- [15] K. Zeng, “Publicly verifiable remote data integrity,” in Proceedings of the 10th International Conference on Information and Communications Security, ser. ICICS '08. Berlin, Heidelberg: Springer-Verlag, 2008, pp. 419-434.
- [16] A. F. Barsoum and M. A. Hasan, “On verifying dynamic multiple data copies over cloud servers,” Cryptology ePrint Archive, Report 2011/447, 2011, 2011, <http://eprint.iacr.org> (Retrieved: 08/06/2013).
- [17] D. Boneh, B. Lynn, and H. Shacham, “Short signatures from the weil pairing,” in ASIACRYPT '01: Proceedings of the 7th International Conference on the Theory and Application of Cryptology and Information Security, London, UK, 2001, pp. 514-532.

**PAPER II****EFFICIENT INTEGRITY VERIFICATION OF REPLICATED  
DATA IN CLOUD USING HOMOMORPHIC ENCRYPTION**

Raghul Mukundan<sup>1</sup>, Sanjay Madria<sup>1</sup>, Mark Linderman<sup>2</sup>

*<sup>1</sup>Department of Computer Science, Missouri University of Science & Technology,  
Rolla, MO 65409*

*<sup>2</sup>Air Force Research Lab, Rome, NY*

**ABSTRACT**

The cloud computing is an emerging model in which computing infrastructure resources are provided as a service over the Internet. Data owners can outsource their data by remotely storing them in the cloud and enjoy on-demand high quality services from a shared pool of configurable computing resources. However, since data owners and the cloud servers are not in the same trusted domain, the outsourced data may be at risk as the cloud server may no longer be fully trusted. Therefore, data integrity is of critical importance in such a scenario. Cloud should let the owners or a trusted third party to check for the integrity of their data storage without demanding a local copy of the data. Owners often replicate their data on the cloud servers across multiple data centers to provide a higher level of scalability, availability, and durability. When the data owners ask the Cloud Service Provider (CSP) to replicate data, they are charged a higher storage fee by the CSP. Therefore, the data owners need to be strongly convinced that the CSP is storing data copies agreed on in the service level contract, and data-updates have been correctly executed on all the remotely stored copies. To deal with such problems, previous multi copy verification schemes either focused on static files or incurred huge update costs in a dynamic file scenario. In this

paper, we propose a Dynamic Multi-Replica Provable Data Possession scheme (DMR-PDP) that prevents the CSP from cheating; for example, by maintaining fewer copies than paid for and/or tampering data. In addition, we also extend the scheme to support a basic file versioning system where only the difference between the original file and the updated file is propagated rather than the propagation of operations for privacy reasons. DMR-PDP also supports efficient dynamic operations like block modification, insertion and deletion on replicas over the cloud servers. Through security analysis and experimental results, we demonstrate that the proposed scheme is secure and performs better than some other related ideas published recently.

## 2.1. INTRODUCTION

When users store data in the cloud, their main concern is whether the cloud can maintain data integrity and whether the data can be recovered when there is data loss or server failure. Cloud Service Providers (CSP), in order to save storage cost, may tend to discard some data or data copies that are not accessed often, or mitigate such data to the second-level storage devices. CSPs may also conceal data loss due to management faults, hardware failures or attacks. Therefore, a critical issue in storing data at untrusted CSPs is periodically verifying whether the storage servers maintain data integrity and store data completely and correctly as stated in the Service Level Agreement (SLA).

Replication is a commonly used technique to increase the data availability in the cloud computing. Cloud replicates the data and stores them strategically on multiple servers located at various geographic locations. Since the replicated data are copies, it is difficult to verify whether the cloud really stores multiple copies of the data. The cloud can easily cheat the owner by storing only one copy of the data. Thus, the owner would like to verify at regular intervals whether the cloud indeed possesses multiple copies of the data as claimed in the SLA. In general, the cloud has the capability to generate multiple replicas when a data

owner challenges the CSP to prove that it possesses multiple copies of the data. Also, it is a valid assumption that the owner of the data may not have a copy of the data stored locally. So, the major task of the owner is not only to verify that the data is intact but also to recover the data if any deletions/corruptions of data are identified. If the owner, during his verification using DMR-PDP scheme, detects some data loss in any of the replicas in the cloud, he can recover the data from other replicas that are stored intact. Since, the replicas are to be stored at diverse geographic locations, it is assumed to be safe that the data loss will not occur at all the replicas at the same time.

Provable data possession (PDP) [2] is a technique to audit and validate the integrity of data stored on remote servers. In a typical PDP model, the data owner generates metadata/tag for a data file to be used later for integrity verification. To ensure security, the data owner encrypts the file and generates tags on the encrypted file. The data owner sends the encrypted file and the tags to the cloud, and deletes the local copy of the file. When the data owner wishes to verify data integrity, he generates a challenge vector and sends it to the cloud. The cloud replies by computing a response on the data and sends it to the verifier/data owner to prove that multiple copies of the data file are stored in the cloud. Different variations of PDP schemes such as [2], [4], [6], [7], [9], [10], [11], [12], [15] were proposed under different cryptographic assumptions. However, most of these schemes deal only with static data files and are valid only for verifying a single copy. A few other schemes such as [3], [5], [8], [13], [14] provide dynamic scalability of a single copy of a data file for various applications which means that the remotely stored data can not only be accessed by the authorized users, but also be updated and scaled by the data owner.

In this paper, we propose a scheme that allows the data owner to securely ensure that the CSP stores multiple replicas. A simple way to make the replicas



look unique and differentiable is by using probabilistic encryption schemes. Probabilistic encryption creates different cipher texts each time the same message is encrypted using the same key. Thus, our scheme uses homomorphic probabilistic encryption to create distinct replicas/copies of the data file and BonehLynnShacham signature scheme (BLS) signatures [17] to create constant amount of metadata for any number of replicas. Probabilistic encryption encrypts all the replicas with the same key. Therefore, in our scheme the data owner will have to share just one decryption key with the authorized users and need not worry about CSP granting access to any of the replicas to the authorized users. The homomorphic property of the encryption scheme helps in efficient file updates. The data owner has to encrypt the difference between the updated file and the old file and send it to the cloud, which updates all the replicas by performing homomorphic addition on the file copies. Any authenticated data structure such as Merkle Hash Trees or Skiplist can be used with our scheme to ensure that the cloud uses the right file blocks for data integrity verification. However, the ways to efficiently manage authenticated data structures in the cloud is not within the scope of this paper. RSA signatures used in [1] can also be used with DMR-PDP scheme for creating data tags. We compare the performance of BLS and RSA signatures and discuss the benefits of using BLS signatures over RSA signatures. We identify the possible attacks the CSP can use to cheat the data owner and provide a security analysis of the proposed protocol against these attacks. Efficiency of the protocol has been experimentally tested and the results demonstrate that our protocol is more efficient than the scheme in [16].

We also extended the DMR-PDP scheme to support a basic file versioning system. The advantage of a file version system is that the data owner can review the changes done to the file and also retrieve the older versions of the file. The cloud computing solutions like Dropbox, Google Drive provide basic file version control system in addition to data storage service. The client-server file versioning models like SVN offer much more features when compared to what these cloud

solutions provide. All these solutions use Delta compression technique where only the difference between the original file and the updated file is propagated and stored on the server. The differences are recorded in discrete files called "deltas" or "diffs". To retrieve any particular version of the file, the change or delta stored on the server is merged with the base version. This reduces the bandwidth and storage space on the server. In our present setting, the data owner uses a homomorphic probabilistic encryption scheme to encrypt and generates multiple encrypted replicas of the file. The cloud stores the replicas across multiple servers. The challenge is to store multiple replicas of the file and maintain the file updates as deltas to support a basic file versioning system when the files are encrypted. To address this, we propose Multiple Replica File Version Control System (MRFVCS) as an extension to the DMR-PDP scheme, which supports encrypted file version control when the data is replicated and the data owner can still use the DMR-PDP scheme to verify the integrity of the data. We implemented the MRFVCS scheme and our experiments show that it is efficient maintaining versions and needs only a little data storage on the data owner side .

Organization: The rest of the paper is organized as follows. Overview of the related work is provided in Section 2 followed by the a detailed description of our scheme in Section 3, and Security analysis in Section 4. MRFVCS is discussed in section 5 followed by Experimental results in Section 6 and Conclusion in Section 7.

## 2.2. RELATED WORK

Ateniese et al. [2] were the first to define the Provable Data Possession (PDP) model for ensuring the possession of files on untrusted storages. They made use of RSA-based homomorphic tags for auditing outsourced data. However, dynamic data storage and multiple replica system are not considered in this scheme. In their subsequent work [3] and [12], they proposed a dynamic version which supports very basic block operations with limited functionality, and does

not support block insertions. In [13], Wang et al. considered dynamic data storage in a distributed scenario, and proposed a challenge-response protocol which can determine data correctness as well as locate possible errors. Similar to [12], only partial support for dynamic data operation is considered. Erway et al. [5] extended the PDP model in [2] to support dynamic updates to stored data files using rank-based authenticated skip lists. However, the efficiency of their scheme remains unclear and these schemes hold good only for verifying a single copy. Wang et al. [14] use Merkle Hash Trees (MHT) for data integrity verification and their scheme supports dynamic data operations, but in their scheme data is not encrypted and works only for a single copy. Hao et al. [8] proposed a scheme that supports both the dynamic data operations and public verifiability. Public verifiability allows anyone, not necessarily the data owner, to verify the integrity of data stored in the cloud by running the challenge-response protocol. Their scheme too do not consider data encryption and their scheme cannot be extended to suit the scenario where multiple data copies are stored.

Curtmola et.al [1] proposed Multiple-Replica PDP (MR-PDP) scheme wherein the data owner can verify that several copies of a file are stored by a storage service provider. In their scheme, distinct replicas are created by first encrypting the data and then masking it with the randomness generated from a Pseudo-Random Function (PRF). The randomized data is then stored across multiple servers. The scheme uses RSA signatures for creation of tags. But, their scheme did not address how the authorized users of the data can access the file copies from the cloud servers noting that the internal operations of the CSP are opaque and do not support dynamic data operations. Ayad F. Barsoum et al. [16] proposed creation of distinct copies by appending replica number to the file blocks and encrypting it using an encryption scheme that has strong diffusion property, e.g., AES. Their scheme supports dynamic data operations but during file updates, the copies in all the servers should be encrypted again and updated on the cloud. This scheme suits perfectly for static multiple replicas but proves costly in a dynamic scenario.

BLS signatures are used for creating tags and authenticated data structures like Merkle Hash Trees are used to ensure that the right file blocks are used during verification. Authorized users of the data should know random numbers in [1] and replica number in [16] to generate the original file.

### 2.3. DYNAMIC MULTI-REPLICA PROVABLE DATA POSSESSION (DMR-PDP) SCHEME

The cloud computing model considered in this work consists of three main components as illustrated in Figure 1: (i) a data owner that can be an individual or an organization originally possessing sensitive data to be stored in the cloud; (ii) a CSP who manages the cloud servers and provides paid storage space on its infrastructure to store the owner's files and (iii) authorized users - a set of owner's clients who have the right to access the remote data and share some keys with the data owner.



Figure 2.1. Cloud Computing Data Storage Model.

**2.3.1. Problem Definition and Design Goals.** More recently, many data owners relieve their burden of local data storage and maintenance by outsourcing their data to a CSP. CSP undertakes the data replication task in order to increase the data availability, durability and reliability but the customers have to pay for using the CSPs storage infrastructure. On the other hand, the cloud

customers should be convinced that the (1) CSP actually possesses all the data copies as agreed upon, (2) integrity of these data copies are maintained, and (3) the owners are able to update the data that they are paying for. Therefore, in this paper, we address the problem of securely and efficiently creating multiple replicas of the data file of the owner to be stored over untrusted CSP and then auditing all these copies to verify their completeness and correctness. Our design goals are summarized below:

1. Dynamic Multi-Replica Provable Data Possession (DMR-PDP) protocols should efficiently and securely provide the owner with strong evidence that the CSP is in possession of all the data copies as agreed upon and that these copies are intact.
2. Allowing the users authorized by the data owner to seamlessly access a file copy from the CSP.
3. Using only a single set of metadata/tags for all the file replicas for verification purposes.
4. Allowing dynamic data operation support to enable the data owner to perform block-level operations on the data files while maintaining the same level of data correctness assurance.
5. Enabling both probabilistic and deterministic verification guarantees.

**2.3.2. Preliminaries and Notations.** In this section, we provide details of the Bilinear mapping and Paillier Encryption schemes used in our present work.

1. Assume that  $F$ , a data file to be outsourced, is composed of a sequence of  $m$  blocks, i.e.,  $F = \{b_1, b_2, \dots, b_m\}$  where  $b_i \in \mathbb{Z}_N$ , where  $\mathbb{Z}_N$  is the set of all residues when divided by  $N$  and  $N$  is a public key in Paillier scheme.
2. Let  $F_i$  represent the file copy  $i$ . So  $F_i = \{b_{i1}, b_{i2}, \dots, b_{im}\}$ , where  $b_{ij}$  represents file block  $b_j$  of file copy  $i$ .

3. BLS Signatures: BLS signatures are short homomorphic signatures that use the properties of bilinear pairings on certain elliptic curves. These signatures allow concurrent data verification, where multiple blocks can be verified at the same time.
4. Bilinear Map/Pairing: Let  $G_1$ ,  $G_2$ , and  $G_T$  be cyclic groups of prime order  $a$ . Let  $u$  and  $v$  be generators of  $G_1$  and  $G_2$ , respectively. A bilinear pairing is a map  $e : G_1 \times G_2 \rightarrow G_T$  with the following properties:
  - Bilinear:  $e(u_1 u_2, v_1) = e(u_1, v_1) \cdot e(u_2, v_1)$ ,  $e(u_1, v_1 v_2) = e(u_1, v_1) \cdot e(u_1, v_2) \forall u_1, u_2 \in G_1$  and  $v_1, v_2 \in G_2$
  - Non-degenerate:  $e(u, v) \neq 1$
  - There exists an efficient algorithm for computing  $e$
  - $e(u_1^x, v_1^y) = e(u_1, v_1)^{xy} \forall u_1 \in G_1; v_1 \in G_2$ , and  $x, y \in \mathbb{Z}_a$
5.  $H(\cdot)$  is a map-to-point hash function:  $\{0, 1\}^* \rightarrow G_1$ .
6. Homomorphic Encryption: A homomorphic encryption scheme has the following properties.
  - $E(m_1 + m_2) = E(m_1) +_h E(m_2)$  where  $+_h$  is a homomorphic addition operation.
  - $E(k * m) = E(m)^k$ .

where  $E(\cdot)$  represents a homomorphic encryption scheme and  $m, m_1, m_2$  are messages that are encrypted and  $k$  is some random number.
7. Paillier Encryption: Paillier cryptosystem is a homomorphic probabilistic encryption scheme. The steps are as follows.
  - Compute  $N = p * q$  and  $\lambda = \text{LCM}(p-1, q-1)$ , where  $p, q$  are two prime numbers.
  - Select a random number  $g$  such that its order is a multiple of  $N$  and  $g \in \mathbb{Z}_N^{2*}$ .

- Public key is  $(N, g)$  and secret key is  $\lambda$ , where  $N = p^*q$ .
- Cipher text for a message  $m$  is computed as  $c = g^m r^N \bmod N^2$  where  $r$  is a random number and  $r \in Z_N^*$ ,  $c \in Z_N^{2*}$  and  $m \in Z_N$ .
- Plain text is obtained by  $m = L(c^\lambda \bmod N^2) * (L(g^\lambda \bmod N^2))^{-1} \bmod N$ .

#### 8. Properties of public key $g$ in Paillier Scheme

- $g \in Z_N^{2*}$ .
- If  $g = (1 + N) \bmod N^2$ , it has few interesting properties
  - (a) Order of the value  $(1 + N)$  is  $N$ .
  - (b)  $(1 + N)^m \equiv (1 + mN) \bmod N^2$ .  $(1 + mN)$  can be used directly instead of calculating  $(1 + N)^m$ . This avoids the costly exponential operation during data encryption.

**2.3.3. DMR-PDP Construction.** In our approach, the data owner creates multiple encrypted replicas and uploads them on to the cloud. The CSP stores them on one or multiple servers located at various geographic locations. The data owner shares the decryption key with a set of authorized users. In order to access the data, an authorized user sends a data request to the CSP and receives a data copy in an encrypted form that can be decrypted using a secret key shared with the owner. The proposed scheme consists of seven algorithms: KeyGen, ReplicaGen, TagGen, Prove, Verify, PrepareUpdate and ExecUpdate. The overview of the communication involved in our scheme is shown in Figure 2.

1.  $(pk, sk) \leftarrow \text{KeyGen}()$ . This algorithm is run by the data owner to generate a public key  $pk$  and a private key  $sk$ . The data owner generates five sets of keys.
  - (a) Key for data tags : This key is used for generating tags for the data. The data owner selects a bilinear map  $e$  and selects a private key  $l \in Z_a$ , where  $l$  is the private key and  $a$  is the order of cyclic groups  $G_1$ ,

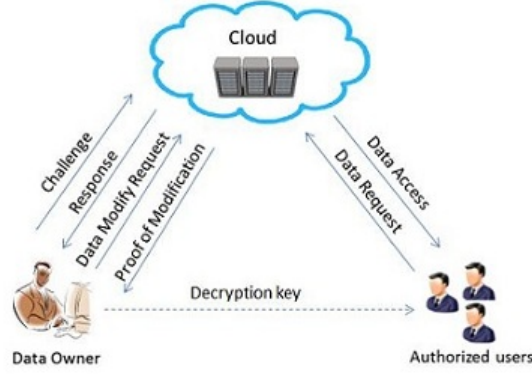


Figure 2.2. DMR-PDP Scheme

$G_2$ , and  $G_T$ . Public key is calculated as  $y = v^j \in G_2$ , where  $v$  is the generator of group  $G_2$ .

- (b) Key for data : This key is used for encrypting the data and thereby creating multiple data copies. The data owner selects paillier public keys  $(N, g)$  with  $g = (1 + N) \bmod N^2$  and secret key  $\lambda$ .
- (c) PRF key for verification: The data owner generates a PRF key  $Key_{PRF}$  which generates  $s$  numbers. These  $s$  numbers are used in creating  $s$  copies of the data. Each number is used in creating one data copy. Let  $\{k_1, k_2, \dots, k_s\} \in \mathbb{Z}_N^*$  be the numbers generated by the PRF key.  $Key_{PRF}$  is maintained confidentially by the data owner and hence the  $s$  numbers used in creating multiple copies are not known to the cloud.
- (d) PRF key for Paillier encryption: The data owner generates a PRF key  $Key_{rand}$ , which is used for generating the random numbers used in Paillier encryption.
- (e) PRF key for Tag generation: The data owner generates a PRF key  $Key_{tag}$ , which is used in generation of tags.

2.  $\{F_i\}_{1 \leq i \leq s} \leftarrow \text{ReplicaGen}(s, F)$ . This algorithm is run by the data owner. It takes the number of replicas  $s$  and the file  $F$  as input and generates  $s$  unique differentiable copies  $\{F_i\}_{1 \leq i \leq s}$ . This algorithm is run by the data owner only once. Unique copies of each file block of file  $F$  is created by encrypting it using a probabilistic encryption scheme, e.g., Paillier encryption scheme.



Through probabilistic encryption, encrypting a file block  $s$  times yields  $s$  distinct cipher texts. For a file  $F = \{b_1, b_2, \dots, b_m\}$  multiple data copies are generated using Paillier encryption scheme as  $F_i = \{(1+N)^{b_1}(k_i r_{i1})^N, (1+N)^{b_2}(k_i r_{i2})^N, \dots, (1+N)^{b_m}(k_i r_{im})^N\}_{1 \leq i \leq m}$ . Using Paillier's properties the above result can be rewritten as  $F_i = \{(1+b_1N)(k_i r_{i1})^N, (1+b_2N)(k_i r_{i2})^N, \dots, (1+b_mN)(k_i r_{im})^N\}_{1 \leq i \leq m}$ , where  $i$  represents the file copy number,  $k_i$  represents the numbers generated from PRF key  $Key_{PRF}$  and  $r_{ij}$  represents random number used in Paillier encryption scheme generated from PRF key  $Key_{rand}$ .  $k_i$  is multiplied by the random number  $r_{ij}$  and the product is used for encryption. The presence of  $k_i$  in a block identifies which copy the file block belongs to. All these file copies yield the original file when decrypted. This allows the users authorized by the data owner to seamlessly access the file copy received from the CSP.

3.  $\phi \leftarrow \text{TagGen}(\text{sk}, F)$ . This algorithm is run by the data owner. It takes the private key  $\text{sk}$  and the file  $F$  as input, and outputs the tags  $\phi$ . We use BLS signature scheme to create tags on the data. BLS signatures are short and homomorphic in nature and allow concurrent data verification, which means multiple data blocks can be verified at the same time. In our scheme, tags are generated on each file block  $b_i$  as  $\phi_i = (H(F) \cdot u^{b_i N + a_i})^l \in G_1$  where  $u \in G_1$ ,  $H(\cdot) \in G_1$  represents hash value which uniquely represents the file  $F$  and  $\{a_i\}_{1 \leq i \leq m}$  are numbers generated from PRF key  $Key_{tag}$  to randomize the data in the tag. Randomization is required to avoid generation of same tags for similar data blocks. The data owner sends the tag set  $\phi = \{\phi_i\}_{1 \leq i \leq m}$  to the cloud.
4.  $P \leftarrow \text{Prove}(F, \phi, \text{challenge})$ . This algorithm is run by the CSP. It takes the replicas of file  $F$ , the tags  $\phi$  and *challenge* vector sent by the data owner as input and returns a proof  $P$  which guarantees that the CSP is actually storing  $s$  copies of the file  $F$  and all these copies are intact. The data owner

uses the proof  $P$  to verify the data integrity. There are two phases in this algorithm:

- (a) **Challenge:** In this phase, the data owner challenges the cloud to verify the integrity of all outsourced copies. There are two types of verification schemes:
- i. Deterministic - here all the file blocks from all the copies are used for verification.
  - ii. Probabilistic - only a few blocks from all the copies are used for verification. A Pseudo Random Function key (PRF) is used to generate random indices ranging between 1 and  $m$ . The file blocks from these indices are used for verification. In each verification a percentage of the file is verified and it accounts for the verification of the entire file.

At each challenge, the data owner chooses the type of verification scheme he wishes to use. If the owner chooses the deterministic verification scheme, he generates one PRF key,  $Key_1$ . If he chooses the probabilistic scheme he generates two PRF keys,  $Key_1$  and  $Key_2$ . PRF keyed with  $Key_1$  generates  $c$  ( $1 \leq c \leq m$ ) random file indices which indicates the file blocks that CSP should use for verification. PRF keyed with  $Key_2$  generates  $s$  random values and the CSP should use each of these random numbers for each file copy while computing the response. The data owner sends the generated keys to the CSP.

- (b) **Response:** This phase is executed by the CSP when a challenge for data integrity verification is received from the data owner. Here, we show the proof for probabilistic verification scheme (the deterministic verification scheme also follows the same procedure). The CSP receives two PRF keys,  $Key_1$  and  $Key_2$  from the data owner. Using  $Key_1$ , CSP generates a set  $\{C\}$  with  $c$  ( $1 \leq c \leq m$ ) random file indices ( $\{C\} \in$

$\{1, 2, \dots, m\}$ ), which indicate the file blocks that CSP should use for verification. Using  $\text{Key}_2$ , CSP generates 's' random values  $T = \{t_1, t_2, \dots, t_s\}$ . The cloud performs two operations; One on the tags and the other on the file blocks.

- i. Operation on the tags: Cloud multiplies the file tags corresponding to the file indices generated by PRF key  $\text{Key}_1$ .

$$\begin{aligned}\sigma &= \prod_{j \in C} (H(F) \cdot u^{b_j N + a_j})^l \\ &= \prod_{j \in C} H(F)^l \cdot \prod_{j \in C} u^{(b_j N + a_j)l} \\ &= H(F)^{cl} \cdot u^{(N \sum_{j \in C} b_j + \sum_{j \in C} a_j)l}\end{aligned}$$

- ii. Operation on the file blocks: The cloud first takes each file copy and multiplies all the file blocks corresponding to the file indices generated by the PRF key  $\text{Key}_1$ . The product of each copy is raised to the power the random number generated for that copy by the PRF key  $\text{Key}_2$ . The result of the above operation for each file copy  $i$  is given by  $(\prod_{j \in C} (1 + N)^{b_j} (k_i r_{ij})^N)^{t_i} \bmod N^2$ . The CSP then multiplies the result of each copy to get the result

$$\begin{aligned}\mu &= \prod_{i=1}^s \left( \prod_{j \in C} (1 + N)^{b_j} (k_i r_{ij})^N \right)^{t_i} \\ &= \prod_{i=1}^s \left( \prod_{j \in C} (1 + N)^{b_j t_i} \prod_{j \in C} (k_i r_{ij})^{N t_i} \right) \\ &= \prod_{i=1}^s \left( (1 + N)^{t_i \sum_{j \in C} b_j} \prod_{j \in C} (k_i r_{ij})^{N t_i} \right) \\ &= \left( \prod_{i=1}^s (1 + N)^{t_i \sum_{j \in C} b_j} \right) \left( \prod_{i=1}^s \left( (k_i)^{c t_i N} \prod_{j \in C} (r_{ij})^{N t_i} \right) \right)\end{aligned}$$

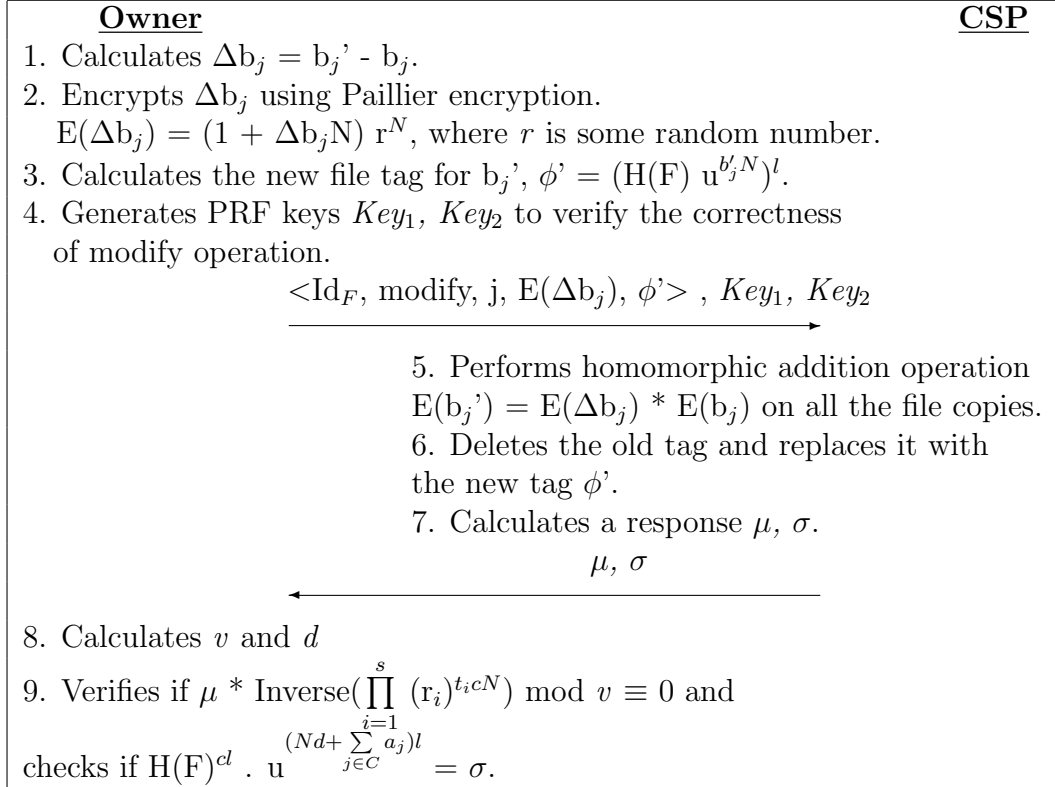


Figure 2.3. Block modification operation in the DMR-PDP scheme

$$= ((1 + N)^{\sum_{i=1}^s t_i \sum_{j \in C} b_j}) (\prod_{i=1}^s (k_i)^{ct_i N}) (\prod_{i=1}^s \prod_{j \in C} (r_{ij})^{N t_i})$$

Using properties of Paillier scheme, the above equation can be rewritten as

$$\mu = (1 + N \sum_{i=1}^s (t_i) \sum_{j \in C} (b_j)) (\prod_{i=1}^s (k_i)^{N c t_i}) (\prod_{i=1}^s (\prod_{j \in C} (r_{ij})^{t_i N}))$$

The CSP sends  $\sigma$  and  $\mu \bmod N^2$  values to the data owner.

5.  $\{1, 0\} \leftarrow \text{Verify}(\text{pk}, P)$ . This algorithm is run by the data owner. It takes as input the public key  $\text{pk}$  and the proof  $P$  returned from the CSP, and outputs 1 if the integrity of all file copies is correctly verified or 0 otherwise. After receiving  $\sigma$  and  $\mu$  values from the CSP, the data owner does the following:

- (a) calculates  $v = (\prod_{i=1}^s (k_i)^{t_i c N})$  and  $d = \text{Decrypt}(\mu) * \text{Inverse}(\sum_{i=1}^s t_i)$ . This can be calculated from the values generated from  $\text{Key}_{rand}$ ,  $\text{Key}_{PRF}$  and the value  $c$ .

- (b) checks if  $\mu * \text{Inverse}\left(\prod_{i=1}^s (r_i)^{t_i c N}\right) \bmod v \equiv 0$ . This ensures that the cloud has used all the file copies while computing the response.
- (c) checks if  $(\text{H}(\text{F})^c \text{u}^{\sum_{j \in C} a_j})^l = \sigma$ . The random numbers in the tag are generated from PRF key  $\text{Key}_{tag}$ . This ensures that the CSP has used all the file blocks while computing the response. If options b and c are satisfied, it indicates that the data stored by the owner in the cloud is intact and the cloud has stored multiple copies of the data as agreed in the service level agreement.

6.  $Update \leftarrow \text{PrepareUpdate}()$ . This algorithm is run by the data owner to perform any operation on the outsourced file copies stored by the remote CSP. The output of this algorithm is an *Update* request. The data owner sends the *Update* request to the cloud that will be of the form  $\langle \text{Id}_F, \text{BlockOp}, j, b_i', \phi' \rangle$ , where  $\text{Id}_F$  is the file identifier,  $\text{BlockOp}$  corresponds to block operation,  $j$  denotes the index of the file block,  $b_i'$  represents the updated file blocks and  $\phi'$  is the updated tag.  $\text{BlockOp}$  can be data modification, insertion or delete operation.

7.  $(\text{F}', \phi') \leftarrow \text{ExecUpdate}(\text{F}, \phi, \text{Update})$ . This algorithm is run by the CSP where the input parameters are the file copies  $\text{F}$ , the tags  $\phi$ , and *Update* request (sent from the owner). It outputs an updated version of all the file copies  $\text{F}'$  along with updated signatures  $\phi'$ . After any block operation, the data owner runs the challenge protocol to ensure that the cloud has executed the operations correctly. The operation in *Update* request can be modifying a file block, inserting a new file block or deleting a file block.

(a) **Modification:** Data modification is one of the most frequently used dynamic operations. The data modification operation in DMR-PDP scheme is shown in Figure 3.

(b) **Insertion:** In the block insertion operation, the owner inserts a new block after position  $j$  in a file. If the file  $F$  had  $m$  blocks initially,

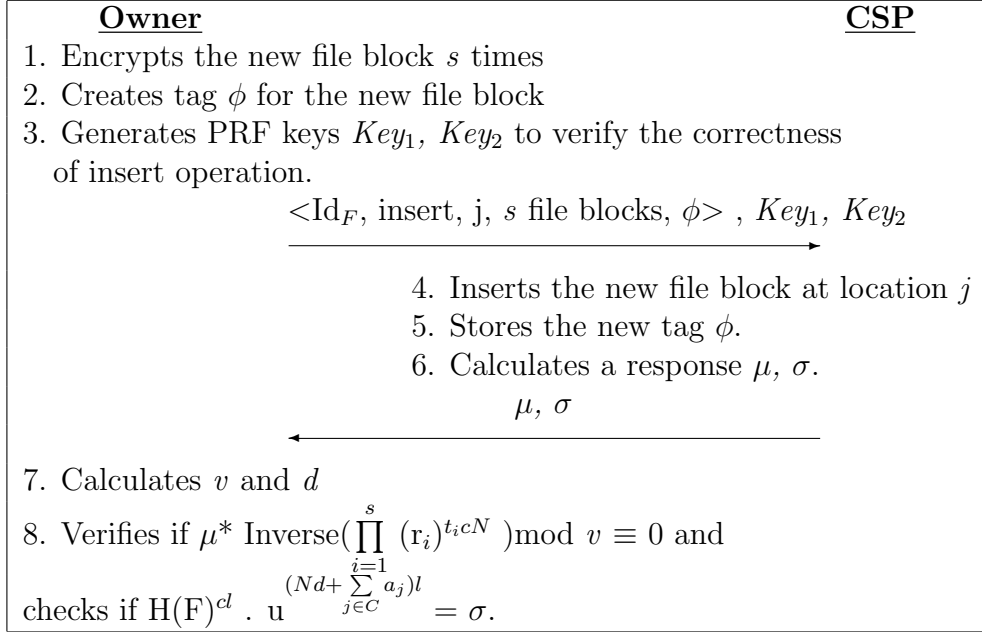


Figure 2.4. Block insertion operation in the DMR-PDP scheme

the file will have  $m+1$  blocks after the insert operation. The file block insertion operation is shown in Figure 4.

- (c) **Deletion:** Block deletion operation is opposite of the insertion operation. When one block is deleted, indices of all subsequent blocks are moved one step forward. To delete a specific data block at position  $j$  from all copies, the owner sends a delete request  $\langle Id_F, \text{delete}, j, \text{null}, \text{null} \rangle$  to the cloud. Upon receiving the request, the cloud deletes the tag and the file block at index  $j$  in all the file copies.

**2.3.4. Using RSA Signatures.** DMR-PDP scheme works well even if RSA signatures are used instead of BLS signatures. The complete DMR-PDP scheme using RSA signatures is shown in Figure 5.

## 2.4. SECURITY ANALYSIS

In this section, we present a formal analysis of the security of our proposed scheme. The data owner encrypts the files and stores them on the cloud which is untrusted and so we identify the cloud as the main adversary in this scheme. The

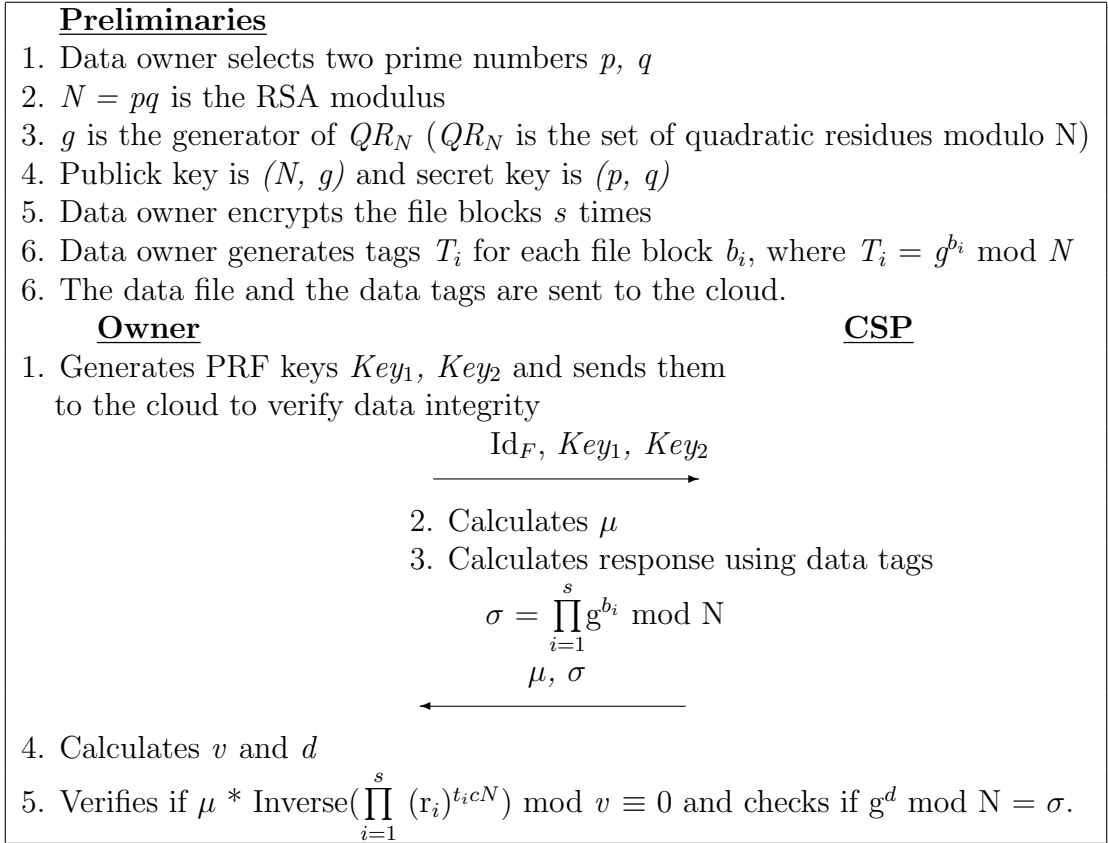


Figure 2.5. DMR-PDP scheme using RSA signatures

scheme is secure only if it does not let the cloud cheat the data owner by deleting file blocks and still pass the challenge/response phase initiated by the data owner.

- **Security against forging the response by the adversary :** In the challenge phase, the data owner sends to the CSP, two PRF keys -  $Key_1, Key_2$  and a parameter 'c' which indicates number of file blocks he wishes to verify. DMR-PDP scheme provides flexibility to the data owner to send different 'c' and PRF keys in each challenge phase to the CSP. This ensures that the response generated by the CSP will not be the same for each challenge sent by the data owner. This eliminates any opportunity for the CSP to forge the response without actually calculating it.
- **Security against deletion of file blocks with same value :** The data tags generated will be the same for similar file blocks. Though the file blocks are encrypted, the cloud can identify similar file blocks by identifying tags with the same value. Cloud can cheat the user by just storing one block and

deleting similar data blocks. To avoid it, DMR-PDP scheme randomizes the data before constructing the tags. The data in the tags are added with random numbers generated from PRF key  $Key_{tag}$ . So, even if the data tag values are the same, the underlying data file block values will not be the same. For file blocks  $b_i = b_j$

$$Tag(b_i) = (H(F).u^{b_i N + a_i})^l$$

$$Tag(b_j) = (H(F).u^{b_j N + a_j})^l$$

where  $a_i, a_j$  are the random numbers generated from PRF key  $Key_{tag}$ . Though the data blocks are the same, the generated tags will differ in value.

## 2.5. MULTIPLE REPLICA FILE VERSION CONTROL SYSTEM (MRFVCS)

MRFVCS is an extension to DMR-PDP scheme to support a basic file versioning system. The data owner encrypts the data, creates multiple replicas and stores them on the cloud. When the data is updated, the data files are not updated directly rather the updates are maintained as deltas. With MRFVCS the data owner can still use the DMR-PDP scheme to verify that the cloud stores multiple replicas and the deltas intact.

**2.5.1. MRFVCS Construction.** The data owner divides the file into multiple file blocks and generates unique multiple replicas and data tags for the file blocks. Unique multiple replicas of the file blocks are generated using the homomorphic probabilistic encryption scheme and the data tags for the file blocks are generated using the BLS signatures as discussed in Section 3.3. The file block replicas and the data tags for the file blocks are sent to the cloud. These encrypted replicas of the file blocks represent the base version of the un-encrypted file blocks. Any modification done to the latest version of un-encrypted file blocks will result in a new version of the file blocks. The new version of the file blocks are not stored



directly on the cloud, and instead, the deltas are stored. Delta is calculated as the difference between the un-encrypted new version of the file block and the un-encrypted base version of the file block. When a particular version of the file blocks are needed, the data owner requests the cloud to merge the delta blocks with the base version of the file blocks to get the required version of the file blocks. The data owner uses a file version table to track the versions of the file blocks. The table is a small data structure stored on the verifier side to validate the integrity and consistency of all files and its versions stored by the CSP. New versions of the file blocks are generated when the data owner performs an update operation on the file blocks. The data update operation includes inserting new file blocks or modifying or deleting a few file blocks. Delta blocks are generated only when the update operation is 'modify'. Once the update operation is done, the file version table is updated by the data owner. The file version table is maintained only on the data owner side. Keeping the file version table only on the data owner side will help the data owner to hide the details of update operations from the CSP. The file version table consists of five columns: Block Number (BN), Delta Block Number (DBN), File Version (FV), Block Version (BV), Block Operation (BO). The BN acts as an indexing to the file blocks. It indicates the physical position of a block in a data file. The DBN is an indexing to the delta block. If delta does not exist, the value is stored as '-'. The FV indicates the version of the entire file and BV indicates the version of the file block. The BO indicates the operation done on the file block. The maximum value of FV gives the latest version of the file and the maximum value of BV for a particular BN gives the latest version of that particular file block. If no entry of a file block number is found in the file version table, it means no update operations are done on the base version of the file block and the file block in the base version and the latest version of the file are same. When a file block with block number B, file version V, and file block version Y is modified, the data owner may choose to change the version of entire file to V+1 or keep the file blocks under the same version. For both these

cases, a new entry is made in the file version table by the data owner. In first case, the table entry will be  $\langle B, -, V+1, 0, \text{Modify} \rangle$  and for the second case, the table entry will be  $\langle B, -, V, Y+1, \text{Modify} \rangle$ . The proposed scheme consists of seven algorithms : Keygen, ReplicaGen, TagGen, Prove, Verify, PrepareUpdate, ExecUpdate, FileVersionRequest, FileVersionDeliver.

1.  $(pk, sk) \leftarrow \text{KeyGen}()$ . Along with the keys described in section 3.3, the data owner generates a PRF key  $Key_{data}$  which is used to randomize the file blocks before encryption.
2.  $\{F'_{B_i}\}_{1 \leq i \leq s} \leftarrow \text{ReplicaGen}(s, F_B)$ . This algorithm is run by the data owner and it slightly differs from the algorithm described in Section 3.3. In this algorithm, the data owner randomizes the data before generating multiple replicas. File blocks are randomized using the random numbers generated from the PRF key  $Key_{data}$ . For a file  $F_B = \{b_1, b_2, \dots, b_m\}$ , the randomized file will be  $F'_B$ , which is  $\{b_1 + x_1, b_2 + x_2, \dots, b_m + x_m\}$  where  $F_B$  represents the base version of the file and  $\{x_j\}_{1 \leq j \leq m}$  are the random numbers generated using  $Key_{data}$ . The data owner uses Paillier encryption, a homomorphic probabilistic encryption scheme, to create  $s$  replicas of the file  $F_B$ . So,  $F'_{B_i} = \{(1+(b_1 + x_1)N)(k_i r_{i1})^N, (1+(b_2 + x_2)N)(k_i r_{i2})^N, \dots, (1+(b_m + x_m)N)(k_i r_{im})^N\}_{1 \leq i \leq s}$ , where  $i$  represents the file copy number,  $k_i$  represents the numbers generated from PRF key  $Key_{PRF}$  and  $r_{ij}$  represents random number used in Paillier encryption scheme generated from PRF key  $Key_{rand}$  (discussed in section 3.3).
3.  $\phi \leftarrow \text{TagGen}(sk, F)$ . This algorithm is run by the data owner and BLS signatures are used to generate the data tags. The details of this algorithm are same as discussed in Section 3.3.
4.  $P \leftarrow \text{Prove}(F_B, F_\Delta, \phi, \text{challenge})$ . This algorithm is run by the CSP. It takes the replicas of file  $F_B$ , all delta files  $F_\Delta$ , tags  $\phi$  and the *challenge* vector sent by the data owner as input and returns a proof  $P$ . Proof  $P$  guarantees

that the CSP is actually storing  $s$  copies of the file  $F_B$  and all the delta files  $F_\Delta$ . The data owner uses proof  $P$  to verify data integrity. The details of this algorithm are the same as discussed in Section 3.3.

5.  $\{1, 0\} \leftarrow \text{Verify}(\text{pk}, P)$ . This algorithm is run by the data owner. It takes as input public key  $\text{pk}$  and the proof  $P$  returned from the CSP, and outputs 1 if the integrity of all file copies is correctly verified or 0 otherwise. The details of this algorithm are same as discussed in Section 3.3.
6.  $Update \leftarrow \text{PrepareUpdate}()$ . This algorithm is run by the data owner to perform any operation on the outsourced file copies stored by the remote CSP. The output of this algorithm is an *Update* request. The data owner sends the *Update* request to the cloud that will be of the form  $\langle \text{Id}_F, \text{BlockOp}, j, b_i', \phi' \rangle$ , where  $\text{Id}_F$  is the file identifier,  $\text{BlockOp}$  corresponds to block operation,  $j$  denotes the index of the file block,  $b_i'$  represents the updated file blocks and  $\phi'$  is the updated tag.  $\text{BlockOp}$  can be data insertion or modification or delete operation.
  - (a) Insertion: An insert operation on any version ' $V$ ' of the file  $F_V$  means inserting new file blocks in the file. The data owner decides the version of the file to which the new file blocks belong, either to current file version  $V$  or to next file version  $V+1$ . If the new file blocks are added to the file version ' $V+1$ ', then ' $V+1$ ' will be the new version of the file. A new entry is made in the file version table as  $\langle \text{BN}, -, V \text{ or } V+1, 0, \text{Insert} \rangle$ . Since there are no delta blocks, the  $\text{DBN}$  value is '-' and since the file blocks are new, the  $\text{BV}$  value is 0.
  - (b) Modification: Modification is done on the latest version of the file blocks. The data owner identifies the block numbers of the file blocks that he wishes to modify and searches the file version table for block numbers. If no entry is found for a particular block number, the file

blocks from the base version are downloaded from the cloud. If an entry is found, then the latest version of the file block is identified and downloaded from the cloud, the file blocks from the base version and the delta blocks associated with the latest version. The downloaded blocks are decrypted and added with the delta to get the latest version of the file blocks. Modify operation is done on the plain text to get the updated plain text. The data owner calculates the new delta as the difference between the updated plain text and the plain text belonging to the base version. Delta is then randomized and then sent to the cloud. Randomization is required in order to not reveal the delta value to the cloud. Let  $M = \{b_i\}$  where  $1 \leq i \leq s$  be the set of file blocks before the update operation and  $M' = \{b'_i\}$  where  $1 \leq i \leq s$  be the file blocks after the update operation. Deltas  $\Delta M$  are calculated as  $\{b'_i - b_i - x_i\}$  where  $1 \leq i \leq s$ . Deltas are randomized using random numbers generated from PRF key  $Key_{data}$ . So,  $\Delta M = \{b'_i - b_i + N - x_i\}$  where  $1 \leq i \leq s$ .  $\Delta M$  values are sent to the cloud.

- (c) Deletion: A delete operation on any version of the file 'V' means deleting few file blocks from the file. The data owner can delete the file blocks from the current version V or delete the file blocks in the next version of the file which is V+1. The data owner makes an entry in the file version table  $\langle BN, -, V \text{ or } V+1, 0, \text{Delete} \rangle$ . The result of the delete operation is just an entry in the file version table while the cloud does not know anything about the delete operation.

7.  $(F', \phi') \leftarrow \text{ExecUpdate}(F, \phi, \text{Update})$ . This algorithm is run by the CSP where the input parameters are the file copies F, the tags  $\phi$ , and *Update* request (sent by the data owner). It outputs new file copies F' along with updated signatures  $\phi'$ . After any block operation, the data owner runs the challenge protocol to ensure that the cloud has executed the operations correctly. The operation in *Update* request can be modifying a file block or

inserting a new file block. The data owner does not send any delete requests to the cloud and so no data blocks will be deleted.

(a) Insertion: New file blocks sent by the data owner are added to the file F.

(b) Modify: Delta values sent by the data owner are stored on the cloud.

8. FileVersionRequest: The data owner identifies the file blocks of the required version by checking the file version table and sends a request to the cloud with the block numbers. The request will be of the form  $\langle \text{BN}, \text{DBN} \rangle$ . The DBN is required to get to the required version of the file. If there is no DBN entry in the file version table, then the request will be  $\langle \text{BN}, - \rangle$ .
9. FileVersionDeliver: For all the file block numbers with DBN value '-' in the FileVersionRequest, the base version of the file is delivered to the data owner. If DBN has a delta file block number, then the cloud encrypts the deltas with the public key and does a homomorphic addition operation on the base version of the file blocks to get the file blocks of the version requested by the data owner. Let  $\Delta M = \{b'_i - b_i + N - x_i\}$  where  $1 \leq i \leq s$ , be the deltas associated with the file blocks of the corresponding file version requested by the data owner. The encrypted deltas  $E(\Delta M)$  is given by  $\{(1 + (b'_i - b_i + N - x_i)N)(r)^N\}$ , where  $r \in \mathbb{Z}_N^*$  is some random number. The cloud then performs a homomorphic addition operation on the requested file blocks on the base version of the file.

$$\begin{aligned}
 E(F_v) &= E(F_b) * E(\Delta M_v). \\
 &= \{(1 + (b_i + x_i)N)(k_i r_i)^N\} * \\
 &\quad \{(1 + (b'_i - b_i + N - x_i)N)(r)^N\}. \\
 &= \{(1 + (b'_i N)(k_i r_i r)^N\}.
 \end{aligned}$$

The file blocks obtained after the homomorphic addition represents the encrypted file blocks of the version requested by the data owner. The encrypted file blocks are sent to the data owner and the data owner decrypts the file blocks to get the version he requested.

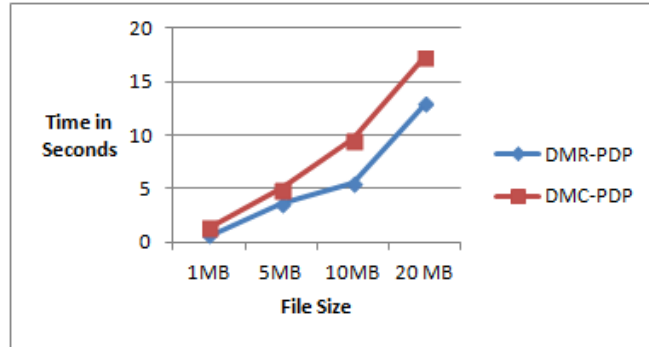
## 2.6. IMPLEMENTATION AND EXPERIMENTAL RESULTS

We implemented our scheme and the protocols in C language. We conducted several experiments using the local cloud servers as well as EC2 cloud instances with different configurations. We measured the computation time for various operations for both the CSP and the user. In addition, we also measured latency in terms of communication cost. Varying file sizes were considered in these experiments.

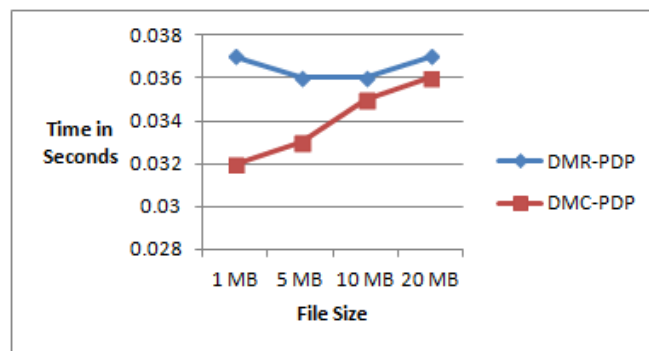
Table 2.1. DMR-PDP Communication cost

Phase	Cost	Data	From	To
Challenge	256 bits	Key <sub>1</sub> , Key <sub>2</sub> , c	Owner	Cloud
Verification	2048 + 160 bits	$\mu, \sigma$	Cloud	Owner
Update	2048 + 160 bits	$b_i', \phi'$	Owner	Cloud

To start with, we conducted several experiments on a system with an Intel(R) Xeon (R) 2.67 GHZ processor and 11 GB RAM running CentOS 6.3. In our implementation, we use PBC library version 0.5.11. To achieve 80-bit security parameter, the curve group with 160-bit group order is selected and the size of modulus N is 1,024 bits. We utilize the Barreto-Naehrig (BN) [18] curve defined over prime field GF(p) with  $|p| = 160$  and embedding degree = 12. The data tags generated are points on this curve and a point on this curve can be represented by 160 bits. We use SHA algorithm for computing file hash, and PBC library provides functions to represent the hash values as a point on the curve. The data owner will have to store three PRF keys of size 128 bit, one secret key for data tags of size 128 bit and one secret key for data encryption of size 1024 bits. The communication cost for each phase incurred in this protocol is shown in Table 1.



(a) CSP computation time

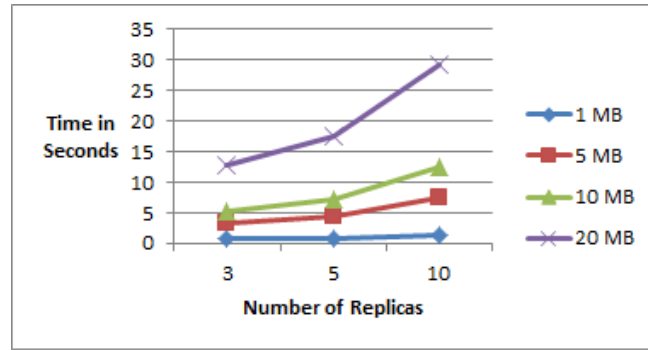


(b) User computation time

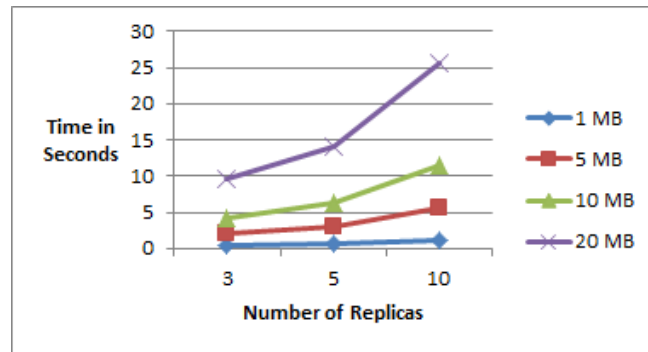
Figure 2.6. Computation time comparison

Here, we compare the performance of the DMR-PDP scheme proposed in this paper with that of the DMC-PDP scheme proposed in [16]. The 1024 bit modulus used for Paillier encryption in this paper is comparable in terms of security to 128 bit AES encryption used in [16]. Figure 6 shows the CSP and User computation times for both the schemes using files of sizes 1, 5, 10 and 20 MB with 3 replicas. The DMR-PDP scheme has lower CSP computation time compared to the DMC-PDP scheme whereas the User computation time for both the schemes differ only in a 1000th of a second. Both the schemes involve just the pairing operation in the user verification phase and hence similar User computation times. Performance of the DMR-PDP scheme is better than that of the DMC-PDP scheme and improves with increase in file size.

The CSP and User computation costs incurred during the response phase of the proposed scheme is depicted in Figure 7.a and 8.a for 1, 5, 10 and 20 MB



(a) CSP computation time for BLS signatures



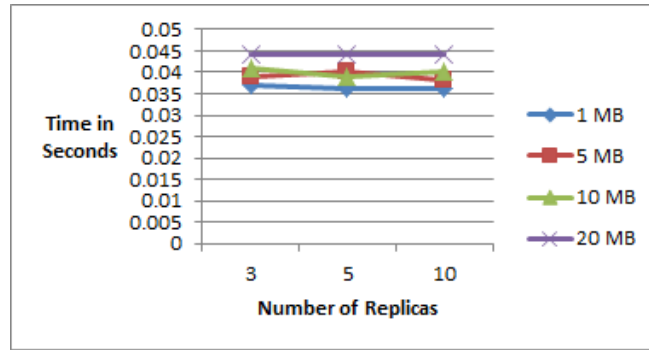
(b) CSP computation time for RSA signatures

Figure 2.7. CSP computation time comparison for the number of replicas on the local cloud servers

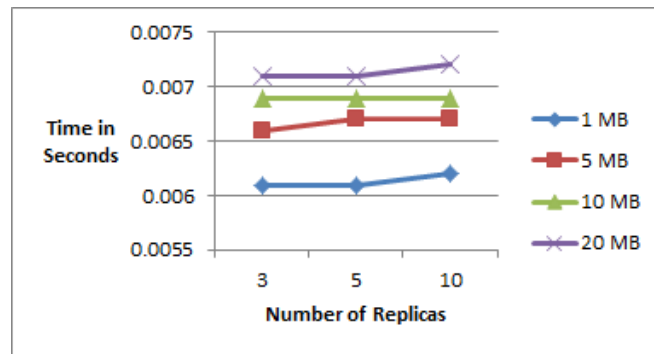
files with 2KB encrypted file block size. Figure 7.a shows the computation time in seconds on the local cloud servers for different number of replicas. For the DMR-PDP scheme the User verification phase involves only one heavy pairing operation and the user computation time is independent of the number of replicas and the file size as shown in Figure 8.a. It has been reported in [1] that if the remote server is missing a fraction of the data, then the number of blocks that needs to be checked in order to detect server misbehavior with high probability is constant and is independent of the total number of file blocks. For example, if the server deletes 1% of the data file, the verifier only needs to check for  $c = 460$  randomly chosen blocks of the file so as to detect this misbehavior with probability larger than 99%. Therefore, in our experiments, we use  $c = 460$  to achieve a high probability of assurance.

Figures 7.b and 8.b show the CSP and the User computation times on the local cloud servers when RSA signatures are used. For better security, we used





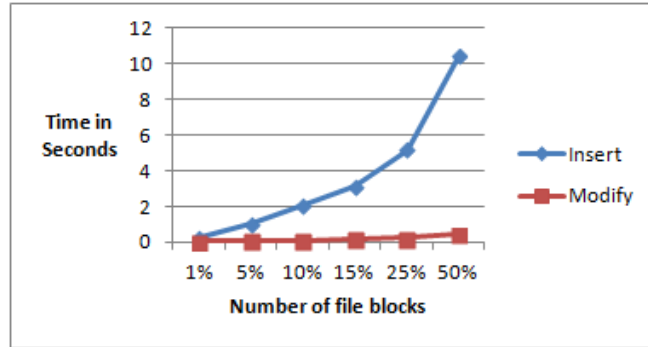
(a) User computation time for BLS signatures



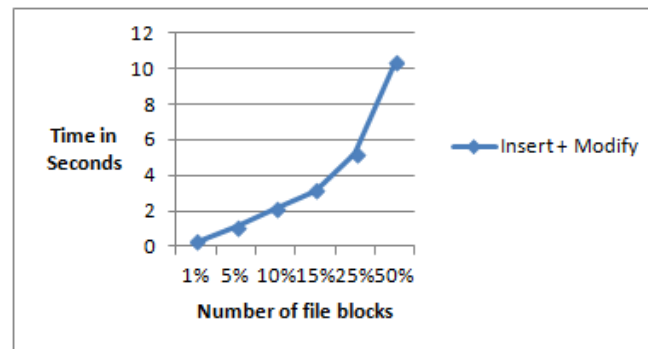
(b) User computation time for RSA signatures

Figure 2.8. User computation time comparison for the number of replicas on the local cloud servers

$N$  to be of size 1024 bits. From Figure 7.a and Figure 7.b, we notice that, the CSP computation time for BLS and RSA signatures is almost the same. Figure 8 shows the comparison of User computation times when BLS and RSA signatures are used. Since user verification of RSA signatures involves only exponential operations and does not involve any complex pairing operations, it is faster than BLS signatures. BLS signatures are better to use with our scheme when compared with RSA signatures since BLS signatures are shorter. The size of RSA signatures is equal to the size of RSA modulus. Since the size of RSA modulus we used is 1024 bits, RSA signatures are also of 1024 bits, whereas size of BLS signatures are just 160 bits. In addition, the BLS construction has the shortest query and response. Further, in our scheme, the data blocks are of size 128 bytes, and so the tag size will also be 128 bytes, if RSA signatures are used. This will increase the communication cost. RSA signatures are useful if the size of the data blocks is huge. [1] uses RSA signatures because they consider data blocks of size 4 KB.



(a) Time for Insert and Modify operations

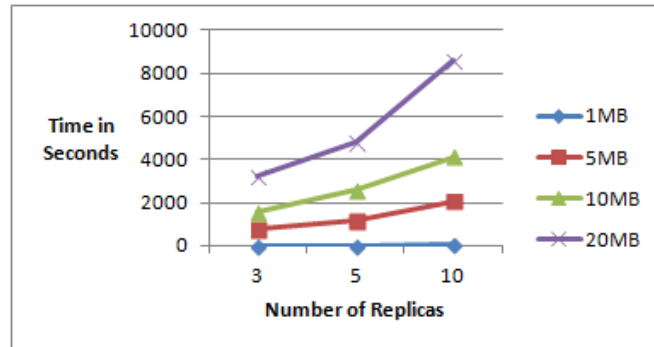


(b) Time for Insert + Modify operations

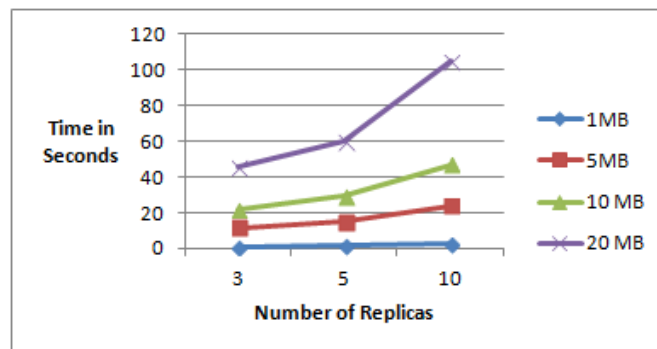
Figure 2.9. Time for file block insert and modify operations on the local cloud servers

The data update operations are done on multiple file blocks at a time by the data owner. The update operation includes file block insert, modify, delete operations in addition to creation of new file tags and their storage on the cloud. We ran the experiments for file block update on a 1 MB file with 3 replicas and file block size of 128 bytes. Figure 9.a. shows the combined User and CSP computation times for file block insert and modify operations, run separately. The experiments are run by inserting and modifying 1% to 50% number of file blocks. For example, a file of size 1 MB has 8192 file blocks. The computation times are calculated by inserting 1% ( $\approx 82$  file blocks) to 50% (4096 file blocks) new file blocks and modifying 1% ( $\approx 82$  file blocks) to 50% (4096 file blocks) of 8192 file blocks. Figure 9.b. shows the computation times when the data owner runs both insert and modify operations on a percentage of file blocks. We notice that modify operations take much less time when compared to insert operations. The time taken for modify operation depends on the time taken for paillier multiplication

of two 256 byte encrypted file blocks whereas the time taken for insert operation depends on the time taken for writing the 256 byte encrypted file blocks on to the hard drive. We do not calculate the time taken for the file block delete operation since the delete operation does not involve any real User and CSP computations.



(a) CSP computation time on EC2 micro instance

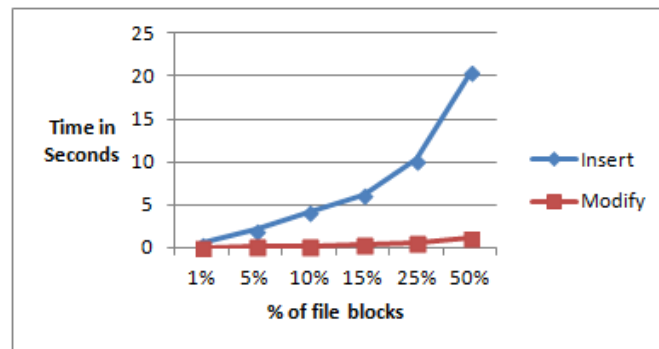


(b) CSP computation time on EC2 large instance

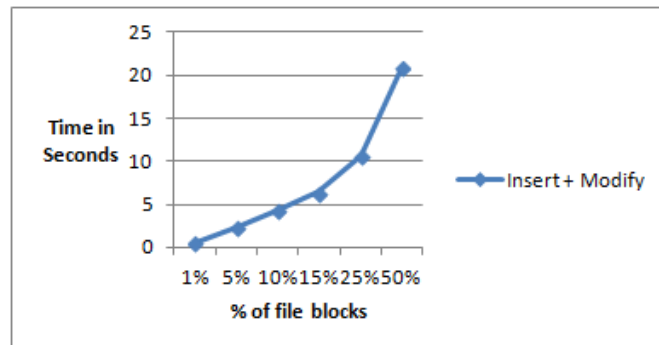
Figure 2.10. CSP computation time for number of replicas on Amazon EC2 instances

We also ran our experiments on a micro and large instance in the Amazon EC2 cloud. We used a 64-bit Ubuntu OS with 25GB storage for the micro and large EC2 instances. A micro instance has 613 MB RAM and uses up to 2 EC2 Compute Units whereas a large instance has 7.5GB RAM and uses 4 EC2 Compute Units. Figure 10 shows the comparison of CSP computation times in micro and large instances on the EC2 cloud. We notice that the experiments run a lot faster on the EC2 large instance when compared to the micro instance. EC2 provides instances which have higher configuration than the large instance and the data owner can use them to get better performance. The user computation time is

independent of the CSP. We calculated it on the local CentOS 6.3 machine which is shown in Figure 8.a. Figures 11 and 12 show the CSP computation times for insert and modify operations on micro and large EC2 instances. It is found that the performance of the update operation is a little faster in large instance when compared to the micro instance. Downloading and uploading file blocks to EC2 micro and large instances take almost the same time. Figure 13 shows the times to download and upload file blocks to EC2 instances.



(a) Time for Insert and Modify operations

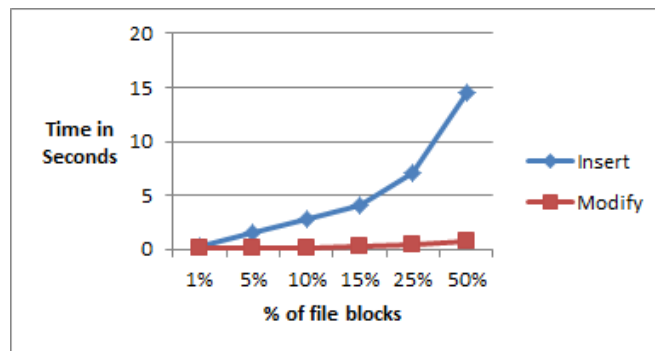


(b) Time for Insert + Modify operations

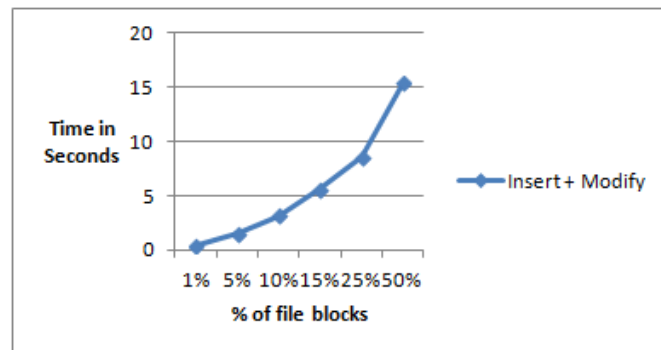
Figure 2.11. Time for file block insert and modify operations in Amazon EC2 micro instance

We also implemented the basic file versioning system 'MRFVCS'. A file version table is created by the data owner to track the data updates. The number of entries in this table depends on the number of dynamic file block operations performed on the data. The file updates are stored as deltas in the cloud. The delta files generated are of size 128 bytes. To get a particular version of the file, the data owner sends 'FileVersionRequest' to the cloud with two parameters <BN,

DBN>. After receiving 'FileVersionRequest', the cloud executes 'FileVersionDeliver' algorithm. For a FileVersionRequest with DBN value '-', the file blocks with block number BN are directly delivered to the data owner and does not involve any CSP computation time. For FileVersionRequest with a valid DBN value, the cloud encrypts the file blocks with block number DBN and does a homomorphic addition operation with file blocks with block number BN. The experiments are run on a 1 MB file on local, Amazon EC2 micro and large instances.



(a) Time for Insert and Modify operations



(b) Time for Insert + Modify operations

Figure 2.12. Time for file block insert and modify operations in Amazon EC2 large instance

Figure 14 shows the time taken by the CSP on various instances for executing 'FileVersionDeliver' algorithm when a number of FileVersionRequests with valid DBN values is sent by the data owner. We considered FileVersionRequests only with valid DBN values since there is no computation time involved for delivering file blocks with DBN value '-'. The number of FileVersionRequests in Figure 14 is represented in terms of percentage of file blocks. For example, a

1 MB file has 8192 file blocks of size 128 bytes. When the data owner sends 81 FileVersionRequests, the CSP encrypts 81 delta blocks (1% of 8192 file blocks) and performs 81 homomorphic addition operations. So any number of FileVersionRequests will lead to the CSP performing operations on those number of file blocks and FileVersionRequests can be represented in terms of number of file blocks. MRFVCS does not involve any computation on the part of the data owner side for executing FileVersionDeliver algorithm and the only cost for the data owner is for maintaining the file version table.

The CSP computation time for executing FileVersionDeliver algorithm in MRFVCS is much less compared to the time taken for update operations in the DMR-PDP scheme.

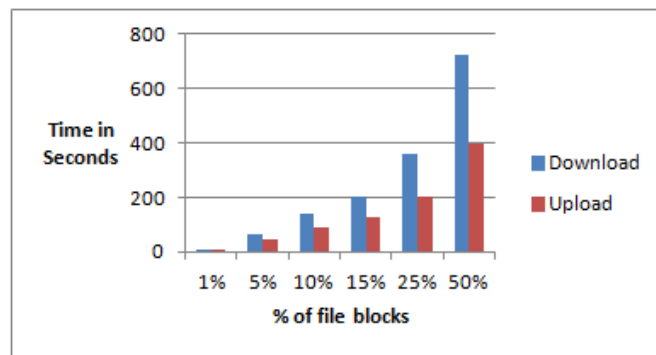


Figure 2.13. Download and Upload time to EC2 Micro and Large Instances

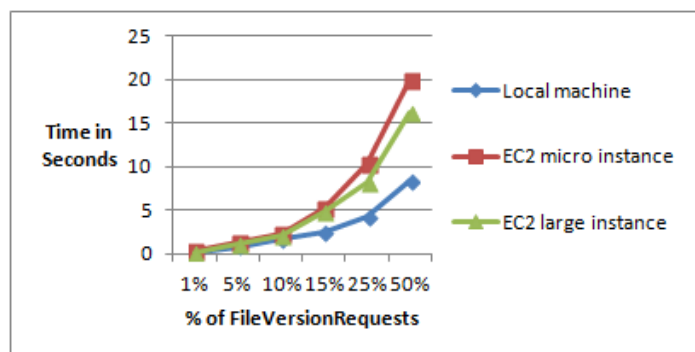


Figure 2.14. CSP computation time comparison for FileVersionDeliver algorithm

## 2.7. CONCLUSION

In this paper, we have presented a scheme for validating the replicated data integrity in a cloud environment. The scheme called Dynamic Multi-Replica Provable Data Possession scheme (DMR-PDP) periodically verifies the correctness and completeness of multiple data copies stored in the cloud. Our scheme considers dynamic data update operations on data copies in the verification process. All the data copies can be decrypted using a single decryption key, thus providing a seamless access to all the data authorized users. The experimental results using the local as well as EC2 cloud instances show that this scheme is better than the previous proposed scheme in terms of dynamic data operations which are performed in much lesser time. In addition, we showed that our scheme works well when extended to support the multiple file versioning where only deltas are stored in the cloud which saves storage cost to the data owner. We believe that these results will help the data owner to negotiate with the cloud provider about the cost and the performance guarantees while maintaining the integrity of the data. Also, these results will provide various incentives to the cloud to take appropriate steps, such as running computations in parallel, to deliver good performance.

**BIBLIOGRAPHY**

- [1] R. Curtmola, O. Khan, R. Burns, and G. Ateniese, “MR-PDP: Multiple-Replica Provable Data Possession,” in 28th IEEE ICDCS, 2008, pp. 411-420.
- [2] G. Ateniese, R. Burns, R. Curtmola, J. Herring, L. Kissner, Z. Peterson, and D. Song, “Provable data possession at untrusted stores,” in CCS '07: Proceedings of the 14th ACM Conference on Computer and Communications Security, New York, NY, USA, 2007, pp. 598-609.
- [3] G. Ateniese, R. D. Pietro, L. V. Mancin, and G. Tsudik, “Scalable and efficient provable data possession,” in SecureComm 08: Proceedings of the 4th International Conference on Security and Privacy in Communication Networks, New York, NY, USA, 2008, pp. 1-10.
- [4] Y. Deswarte, J.-J. Quisquater, and A. Sadane, “Remote integrity checking,” in 6th Working Conference on Integrity and Internal Control in Information Systems (IICIS), S. J. L. Strous, Ed., 2003, pp. 1-11.
- [5] C. Erway, A. Kupcu, C. Papamanthou, and R. Tamassia, “Dynamic provable data possession,” in CCS 09: Proceedings of the 16th ACM Conference on Computer and Communications Security, New York, NY, USA, 2009, pp. 213-222.
- [6] D. L. G. Filho and P. S. L. M. Barreto, “Demonstrating data possession and uncheatable data transfer,” Cryptology ePrint Archive, Report 2006/150, 2006 (Retrieved: 08/06/2013).
- [7] P. Golle, S. Jarecki, and I. Mironov, “Cryptographic primitives enforcing communication and storage complexity,” in FC'02: Proceedings of the 6th International Conference on Financial Cryptography, Berlin, Heidelberg, 2003, pp. 120-135.
- [8] Z. Hao, S. Zhong, and N. Yu, “A privacy-preserving remote data integrity checking protocol with data dynamics and public verifiability,” IEEE Transactions on Knowledge and Data Engineering, vol. 99, no. PrePrints, 2011.
- [9] E. Mykletun, M. Narasimha, and G. Tsudik, “Authentication and integrity in outsourced databases,” Trans. Storage, vol. 2, no. 2, 2006.
- [10] F. Sebe, J. Domingo-Ferrer, A. Martinez-Balleste, Y. Deswarte, and J.-J. Quisquater, “Efficient remote data possession checking in critical information infrastructures,” IEEE Trans. on Knowl. and Data Eng., vol. 20, no. 8, 2008.
- [11] M. A. Shah, M. Baker, J. C. Mogul, and R. Swaminathan, “Auditing to keep online storage services honest,” in HOTOS'07: Proceedings of the 11th USENIX workshop on Hot topics in operating systems, Berkeley, CA, USA, 2007, pp. 1-6.



- [12] M. A. Shah, R. Swaminathan, and M. Baker, "Privacy-preserving audit and extraction of digital contents," Cryptology ePrint Archive, Report 2008/186, 2008.
- [13] C. Wang, Q. Wang, K. Ren, and W. Lou, "Ensuring data storage security in cloud computing," Cryptology ePrint Archive, Report 2009/081, 2009, <http://eprint.iacr.org> (Retrieved: 08/06/2013).
- [14] Q. Wang, C. Wang, J. Li, K. Ren, and W. Lou, "Enabling public verifiability and data dynamics for storage security in cloud computing," in ESORICS09: Proceedings of the 14th European Conference on Research in Computer Security, Berlin, Heidelberg, 2009, pp. 355-370.
- [15] K. Zeng, "Publicly verifiable remote data integrity," in Proceedings of the 10th International Conference on Information and Communications Security, ser. ICICS '08. Berlin, Heidelberg: Springer-Verlag, 2008, pp. 419-434.
- [16] A. F. Barsoum and M. A. Hasan, "On verifying dynamic multiple data copies over cloud servers," Cryptology ePrint Archive, Report 2011/447, 2011, 2011, <http://eprint.iacr.org> (Retrieved: 08/06/2013).
- [17] D. Boneh, B. Lynn, and H. Shacham, "Short signatures from the weil pairing," in ASIACRYPT '01: Proceedings of the 7th International Conference on the Theory and Application of Cryptology and Information Security, London, UK, 2001, pp. 514-532.
- [18] P. S. L. M. Barreto and M. Naehrig, "Pairing-friendly elliptic curves of prime order", in Proceedings of SAC 2005, volume 3897 of LNCS. Springer-Verlag, 2005, pp. 319-331.
- [19] Raghul Mukundan, Sanjay Madria and Mark Linderman, "Replicated data integrity verification in cloud ", in IEEE Data Engineering Bulletin 2012.

## SECTION

### 3. CONCLUSION

In this work, we have presented a scheme for validating the replicated data integrity in a cloud environment. The scheme called Dynamic Multi-Replica Provable Data Possession scheme (DMR-PDP) periodically verifies the correctness and completeness of multiple data copies stored in the cloud. Our scheme considers dynamic data update operations on data copies in the verification process. All the data copies can be decrypted using a single decryption key, thus providing a seamless access to all the data authorized users. The experimental results using the local as well as EC2 cloud instances show that this scheme is better than the previous proposed scheme in terms of dynamic data operations which are performed in much lesser time. In addition, we showed that our scheme works well when extended to support the multiple file versioning where only deltas are stored in the cloud which saves storage cost to the data owner. We believe that these results will help the data owner to negotiate with the cloud provider about the cost and the performance guarantees while maintaining the integrity of the data. Also, these results will provide various incentives to the cloud to take appropriate steps, such as running computations in parallel, to deliver good performance.

**BIBLIOGRAPHY**

- [1] Provable Data Possession at Untrusted Stores - Giuseppe Ateniese, Randal Burns, Reza Curtmola, Joseph Herring, Lea Kissner, Zachary Peterson Dawn Song, ACM Conference on Computer and Communications Security (CCS) 2007
- [2] Privacy-Preserving Public Auditing for Secure Cloud Storage - Cong Wang, Sherman S.-M. Chow, Qian Wang, Kui Ren, and Wenjing Lou, The 29th IEEE Conference on Computer Communications (INFOCOM'10), San Diego, CA, March 15-19, 2010.
- [3] Ensuring Data Storage Security in Cloud Computing - Cong Wang, Qian Wang, and Kui Ren, Wenjing Lou, The 17th IEEE International Workshop on Quality of Service (IWQoS'09), Charleston, South Carolina, July 13-15, 2009.
- [4] Proofs of Retrievability: Theory and Implementation - Kevin D. Bowers, Ari Juels, Alina Oprea, In: Proc. of ACM-CCSW '09, pp.43-54, 2009.
- [5] Provable Possession and Replication of Data over Cloud Servers - Ayad F.Barsoum and M.Anwar Hasan, Centre For Applied Cryptographic Research (CACR), University of Waterloo, Report 2010/32, 2010, <http://www.cacr.math.uwaterloo.ca/techreports/2010/cacr2010-32.pdf> (Retrieved: 08/06/2013).
- [6] R. Curtmola, O. Khan, R. Burns, and G. Ateniese, MR-PDP: Multiple-Replica Provable Data Possession, in 28th IEEE ICDCS, 2008, pp. 411420.
- [7] C. Erway, A. Kup, C. Papamanthou, and R. Tamassia, Dynamic provable data possession, in CCS 09: Proceedings of the 16th ACM Conference on Computer and Communications Security, New York, NY, USA, 2009, pp. 213222.
- [8] A. F. Barsoum and M. A. Hasan, On verifying dynamic multiple data copies over cloud servers, Cryptology ePrint Archive, Report 2011/447, 2011, 2011, <http://eprint.iacr.org> (Retrieved: 08/06/2013).

## VITA

Raghul Mukundan was born on February 19, 1987 in Chittor, India. He received distinction in Bachelor of Technology degree in Computer Science and Engineering from Amrita Vishwavidyapeetham, India in 2008. He has been a graduate student in the Computer Science Department at Missouri University of Science and Technology since August 2011 and worked as a Graduate Research Assistant under Dr. Sanjay Kumar Madria from August 2011 to July 2013. He received his Masters in Computer Science at Missouri University of Science and Technology in August 2013.

