# Secure Data Storage Optimization Over Cloud using Encrypted Cloud Data Deduplication Technique

**Dr. Harsh Lohia, Suhas A. Lakade, Mr. Yuvraj R. Gurav**

Department of Computer Science and Engineering, School Of Engineering , Sri Satya Sai University Of Technology And Medical Science, Sehore (MP), India Dr. Harsh Lohiya e-mail : lohiya27harsh@gmail.com, , Suhas Lakade, e-mail: suhaslakade@gmail.com, Yuvraj Gurav e-mail : yuvi1333@gmail.com

*To Cite this Article*

**Dr. Harsh Lohia, Suhas A. Lakade, Mr. Yuvraj R. Gurav, "Secure Data Storage Optimization Over Cloud using Encrypted Cloud Data Deduplication Technique"** *Journal of Science and Technology, Vol. 09, Issue 01,- Jan 2024, pp131-138*

**Abstract**

Data deduplication is a crucial data compression technique used to eliminate duplicate copies of repetitive data, commonly employed in cloud storage to reduce storage space and bandwidth usage. To ensure the confidentiality of sensitive data while supporting deduplication, a method known as convergent encryption has been introduced to encrypt data before outsourcing. In this paper, we present the first formal approach to address the problem of authorized data deduplication. Unlike traditional deduplication systems, we also consider the varying privileges of users during the duplicate check process, in addition to the data itself. We introduce several novel deduplication constructions that facilitate authorized duplicate check within a hybrid cloud architecture. A comprehensive security analysis validates the effectiveness of proposed scheme based on the specified security model. As a proof of concept, we have implemented a prototype of the proposed authorized duplicate check scheme and performed testbed experiments using it. The results demonstrate that proposed approach incurs minimal overhead compared to normal operations.

## 1 Introduction

The Businesses and individuals are encouraged to outsource their data to a cloud server due to the considerable flexibility and cost advantages that come with cloud computing. IDC's analytical report predicts that by 2020, there will be more data than 44 ZB globally. The need for novel methods to effectively use storage space and network bandwidth is critical given the rapid expansion of data quantities stored in cloud servers. Data deduplication has drawn a lot of interest from both academics and business in order to accomplish this. For instance, client-side deduplication is a technique used by Google Drive, SpiderOak, and Dropbox to lower storage costs and conserve network traffic. Deduplication techniques, on the other hand, use data similarity (at the file or block level) to identify duplicate data and save storage space by storing just one copy of the data in the cloud server. However, from a security standpoint, data deduplication faces new difficulties due to cross-user sharing of data among data owners.

The data owners decide to encrypt the data with their own keys before outsourcing it to the cloud due to privacy and confidentiality concerns. As a result, cross-user deduplication is difficult since the same data encrypted with multiple data owners' keys will produce different ciphertexts, making it impossible for the cloud to detect more duplication. Convergent Encryption (CE), which use a hash value of the data as the encryption key, is suggested as a solution to this issue. As a result, the same data will be encrypted into the same ciphertext, allowing for ciphertext deduplication. Even though it first appears

to be a viable option for achieving both confidentiality and deduplication at the same time, it is regrettably plagued by well-known flaws such tag consistency issues and brute-force attacks. Message-locked encryption (MLE) and leakage-resilient deduplication techniques are suggested as solutions to these issues.

Incross-user deduplication, dynamic ownership management is being studied further. Owners of the data may ask the cloud server to erase or amend their data as time goes on. The owners of the revoked data should not be able to access the associated data kept in the cloud storage after the cloud server processes these requests (forward secrecy). In addition, when a data owner uploads data that has already been stored in a cloud service, the data owner should be permitted access to the data after acquiring ownership (backward secrecy). [2] suggest their dynamic ownership management technique for cross-user file-level deduplication in order to satisfy these security criteria by re-encrypting the ciphertext. These dynamic ownership changes could, however, take place often and add a lot of computational cost. Additionally, the encryption key is generated deterministically and hardly ever modified after the initial key generation in CE-based data deduplication systems. Therefore, a research challenge is how to effectively and securely establish ownership management in cross-user data deduplication.

Since the cloud server might not be able to confirm consistency between the ciphertext and the corresponding plaintext for user-side randomized convergent encryption, the valid data owner should download the relevant ciphertext to calculate the randomized convergent key and to check the tag (RCE). The ciphertext must be uploaded and downloaded, and there is significant computational cost associated with the decryption and re-encryption operations. A maliciously created ciphertext replacement attack or a duplicate fake assault may also be launched by an untrustworthy data owner who only has access to the data's tag in order to persuade the cloud server that they are the true owners of the whole encrypted material.

Bloom filter-based PoW strategy, which is more effective at the user side and at the size of cloud servers, may be a viable option to analyze this problem. However, under the dynamic ownership management system, the owner of revoked data may preserve all legitimate tag proofs in order to regain possession without the necessary file. We should redesign the PoW to support update, which is not achieved or even taken into consideration in existing methods, in order to resist this rejoin assault 1. As a result, another difficulty is how to create a mutual PoW verification system that is secure and effective while supporting ownership management through ciphertext deduplication.

## 2 Review Of Literature

This section describes concisely the detail methodology/procedures employed so that anyone wishing to replicate the trial can do so and obtain comparable results. Provide sufficient details as to remove any possible ambiguities with respect to design, treatments, measurements, analysis, etc. Where methods employed are commonly known in a given field details should be omitted and the reference given instead. Modifications to known methodology must however be clearly described and explained. One paragraph must not contain only one sentence.

### 2.1 Literature Survey

To deduplicate encrypted data stored in the cloud, give a representative re-encryption and ownership challenge. Users can communicate data even while they are offline because to these solutions, which is their main benefit. The key drawback of these methods is that CSP must be properly designed in order for deduplication management to work.[1]

This system is made available, which offers safe deduplicated storage to fend off brute-force attacks. In DupLESS, a Key Server that is distinct from a Storage Service helps a group of linked clients encrypt their data. The fundamental benefit of these methods is the avoidance of brute force assaults and the ability for clients to encrypt data using the key server, which differs from a separate storage server. The fundamental drawback of these strategies is that they cannot offer flexibility to different data consumers.[2]

A policy-based deduplication representative scheme was put forth, however it did not assess the performance of the scheme or take into account the deletion and owner management examples of duplicated data management. This method's key benefit is that it uses policy-based deduplication to authorize trust relationships between cloud storage components. The fundamental disadvantage of this method is that policy-based deduplication does not take data deletion and owner management into account.[3]

As far as file storage is concerned, data deduplication, optimum node selection, server load balancing, and server load optimization, provision index name servers to carry out file compression, chunk matching, real-time feedback control, IP information, and busy level index monitoring. This method's primary resource is the Index Name Servers algorithm's advice for lowering resource demands and enhancing system performance. Additionally, INS manages server load balancing. This method's main drawback is that encrypted data cannot be deduplicated.[4]

Under the assumption that CSP is aware of the data's encryption key, they have endorsed the system. Since the CSP cannot be fully trusted by the owners or holders of the data, it cannot be used in that circumstance. The primary benefit of these techniques is that they support both plaintext and ciphertext deduplication. The primary drawback of these methods is that they don't offer encrypted data deduplication.[5]

To verify the integrity of the distant data, introduce a formal PoR model. Their plan makes use of error-correcting coding and cloak blocks called sentinels that are concealed within the original file blocks to detect data damage and guarantee both custody and retrievability of files at distant cloud servers.[6]

Establish that, regardless of the firewall settings of the compromised machine, deduplication systems could be used as a covert conduit for malicious software to connect with the outside world and as a side channel to provide sensitive information about the contents of files. They suggest employing encryption to solve the issue, however since this will virtually prohibit deduplication from occurring, the encryption may not be acceptable. By establishing a random threshold for every file and only doing deduplication if the number of copies of the file exceeds this threshold, they hope to reduce the security risks of deduplication. But rather of eliminating the security issues, their method just lessens them.[7]

**2.2 Summary of Literature Survey**

The proof-of-ownership concept, in which a client can prove to a server that it genuinely owns a copy of a file without actually uploading it by using the Merkle tree and error-control encoding. Their plan, nevertheless, cannot ensure the validity of the evidence in each challenge. Furthermore, their plan necessitates the wasteful construction of the Merkle tree using the encoded data. How to check the integrity of memory contents in computing platforms and make sure that memory corruption doesn't impair the computations being done is another closely connected subject that has been extensively investigated. The Merkle tree is often used in these techniques to authenticate memory and data.

**3 Proposed System**

**3.1 Motivation**

The cloud computing paradigm is the information technology industry's next-generation architecture, and it offers consumers significant savings on compute, storage, bandwidth, and transmission expenses. To significantly reduce costs for the CSP, cloud technology moves all data, databases, and software over the Internet. There are currently services available in cloud computing such as Infrastructure-as-a-Service (IaaS), Platform-as-a-Service (PaaS), Software-as-a-Service (SaaS), and Database-as-a-Service (DaaS). The DaaS attribute of cloud computing, namely cloud data storage services, is the focus of this argument. Popular cloud storage services like SpiderOak, Mozy, and Dropbox are a few examples The growth of digital content has become unstoppable at both the corporate and individual levels with the introduction of cloud computing and associated digital storage services. Additionally, the usage of the Internet and other digital services has resulted in an explosion of digital data, particularly that stored in cloud storage. It would be advantageous in terms of storage savings if the replicated data could be deleted from the data storages given the enormous challenge of big data. Data privacy concerns related to the cloud paradigm are a significant barrier preventing user acceptance of cloud storage services. Data deduplication is a useful method used to solve the problem. The main need is that the data deduplication should be built to meet the system under consideration's requirements for efficiency and security.

**3.2 Aims and Objectives**

The research goals of our work are as follows in order to establish an efficient, secure cloud data deduplication system with ownership management:

- **Aim:**

To provide a secure, effective, and ownership management cloud data deduplication solution

- **Objectives:**

1. To facilitate data confidentiality and dynamic ownership management, we offer an effective and secure data deduplication system.

2. To resist the poison attack and the duplicate faking attack of cross-user file-level deduplication while using the network bandwidth effectively during the data deduplication, we build a novel PoW strategy based on Bloom filter.

3. To lower the update frequency and computational burden for cross-user file level deduplication while ensuring the security of the uploaded data, we develop a lazy update technique.
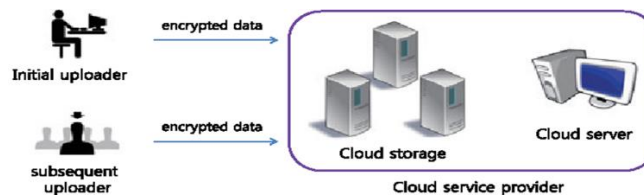
### 3.3 Methodology



**Fig. 1** Caption The system model for our data deduplication system.

The client (user) who owns the data and subcontracts it to the cloud service provider is the data owner. A data owner's credentials and passwords should be validated before granting access to a cloud service provider. The data is then encrypted and uploaded to the cloud service provider along with the appropriate index information, also known as a tag. A subsequent uploader is referred to if the uploaded data already exist in the cloud service provider; otherwise, an initial uploader is used. In our system, after passing the mutual PoW verification, the succeeding uploader does not upload any duplicate material in order to conserve network traffic.

A company that offers cloud storage services is the CSP. It is made up of a cloud storage and a cloud server. If necessary, the cloud server deduplicates the data that has been outsourced from the data owners and saves it in the cloud storage. In order to achieve access control, the cloud server maintains a particular type of data structure with a password and credential for each data owner. In addition, the cloud server also keeps ownership lists, which are made up of a tag for the data that is stored and the names of the owners of that data. The cloud server handles the dynamic ownership based on the lists and restricts access to the data that is stored. Here, we develop a lazy update technique to lessen the frequency of cross-user update operations.

Be aware that our system supports both inside-user block level deduplication and cross-user file level deduplication. The data owner must first execute a cross-user file level duplicate check before uploading a file. In the event of a duplicate file, the data owner performs file-level deduplication; in the absence of a duplicate file, the data owner may also be selected as the initial uploader or may perform inside-user block-level deduplication and determine which unique blocks need to be uploaded.

- **Security Requirements:**

The plaintext of the ciphertext saved in the cloud storage shouldn't be accessible to unauthorized data owners who can't substantiate ownerships.

In order to protect against the duplicate faking attack and the poison attack, the deduplication method should ensure tag consistency. That example, the deduplication technique should offer protection against duplicate faking attacks, in which a genuine file is covertly replaced by a fake one bearing the same tag, enabling an attacker to decrypt the data even if they only have access to the tag.

The most crucial statistic for data deduplication, aside from security concerns, is efficiency. Since there is a lot of data, it is important to maximize the bandwidth for outsourcing data to the cloud server, the storage efficiency, the transmission efficiency, and the computing efficiency (the overhead for dynamic ownership management).

### 3.3.1 System Modelling

**1. User registration and login**

The user can allow to login when the group has generated a team to validate the user and process him towards the group. In this module, a new user can be permitted to register by supplying the user details like name, email, age, etc. It will be taken care of to validate user information and add a new entry to the registration database.

**2. User joining the group and file upload**

Each user will have a unique key that they can use to join the group after being authenticated. In the case of a file upload, the user selects a file from his computer and creates a hash key for each file. The production of hash keys is accommodated to prevent file duplication to the cloud. The user cannot upload a file that is already in the cloud.

**3. File encryption and storage in cloud**

After the cloud has verified the user's file, we employ a cryptographic approach to raise the level of security there. We use the ECC algorithm, which changes a file to binary, encrypts it, and stores it in the cloud. The information that is kept in the cloud will be encrypted.

**4. User file request and download**

Any user who has already signed up and joined the group with a working key may ask for the file to be uploaded to the cloud. After verifying the user's identity, the cloud service provider can accept a file request, decrypt it using the ECC algorithm, and then transfer the requested file to the user. The file will then be downloaded to the user's computer.

### 3.4 Expected Outcomes

The expected outcome of the deduplication mechanism on cloud are as follows:

1. **Improved Storage Efficiency:** Data deduplication eliminates redundant data across different files, ensuring that only unique data is stored in the cloud. This results in significant storage space savings and efficient utilization of cloud resources. As a consequence, organizations can store more data while minimizing storage costs.

2. **Reduced Bandwidth Consumption:** Deduplication reduces the amount of data transferred to and from the cloud, lowering the bandwidth requirements. This is particularly important for organizations with limited network resources or those paying for data transfer.

3. **Enhanced Data Privacy and Security:** Encrypting data before deduplication ensures that sensitive information remains confidential even when it is stored on the cloud. Deduplication techniques that work on encrypted data ensure that identical encrypted blocks are treated as duplicates, preserving data privacy.

4. **Faster Data Backup and Recovery:** With deduplication in place, backups are faster as only new or unique data needs to be transferred. In case of data loss, recovery times are shorter since there's less data to retrieve.

5. **Cost Savings:** By optimizing storage and reducing bandwidth usage, organizations can save money on cloud storage and data transfer costs. This is especially relevant for businesses with large volumes of data.

6. **Scalability:** The technique allows for efficient scalability as it optimizes data storage. Organizations can easily accommodate growing data volumes without incurring exorbitant costs.

### 4 Results

The expected vs actual results of the proposed system are described below:

| Requirement | Expected Output | Actual Output |
|---|---|---|
|  |  |  |

| | | |
|---|---|---|
| Cloud data must be checked for duplicate chunks if any | Proposed system should deny user from uploading identical duplicate data (chunks) over cloud. | The proposed system divides the file being uploaded into chunks based on sentences. It verifies the availability of duplicate chunks uploaded over the cloud. |
| Optimal storage space allocation | Data should be securely encrypted and stored over the cloud. | The system maintains numerical calculation of the percentage of duplicate data uploaded and plots it on a graph. |
| Secure encryption of stored data | Non-authorized users must not be allowed to login to the system. | The system enforces secure encryption of stored data to protect it from unauthorized access. |
| Access restricted to users providing valid ownership proof | Only the user providing valid ownership proof should access the data. | The system implements access control measures to ensure that only users with valid ownership proof can access the data. |
| Graphical representation of duplicate data analysis | The graph shows the percentage of duplicate contents from the uploaded files. | The graph shows the percentage of duplicate data uploaded, based on the analysis of the file chunks. |

The proposed system aims to ensure efficient storage utilization and secure data management in the cloud. To achieve this, it checks for duplicate chunks of data before storing them and denies users from uploading identical duplicate data. This helps to minimize storage redundancy and optimize space allocation. Additionally, the system ensures the secure encryption and storage of data, protecting it from unauthorized access.

Access control measures are implemented to restrict system login to only authorized users who can provide valid ownership proof. This ensures that sensitive data can only be accessed by individuals with the appropriate permissions. The system also provides a graphical representation of the analysis, showing the percentage of duplicate data uploaded.. The analysis reveals that 66% of the data is found to be duplicate, while only 33% of unique data chunks are uploaded.
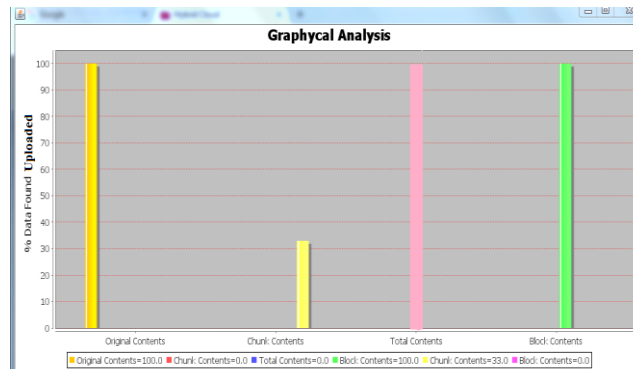
**Fig.02**.Graphical ResultAnalysis.

Overall, the proposed system prevents user from uploading duplicated on cloud. Data stored on cloud must be in secure encrypted format. Malicious user not able to upload or download data on cloud. The user who has proof of ownership only that user can modify data.



**Fig. 03** Home

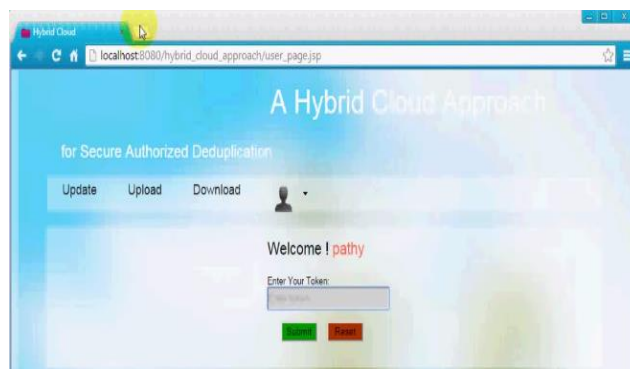The results in Figure 3 depicts the landing page of the proposed system.



**Fig.04** UserWindow

The result page in Figure 4 depicts the page where user enters his/her identification token allocated to the user while registering on the cloud to upload the files on the cloud.
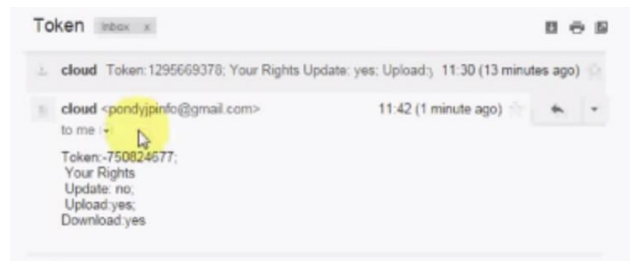
**Fig. 5**. ActivationToken Mail

The result in Figure5 depicts that the activation token which is sent to the user over the email with corresponding details.
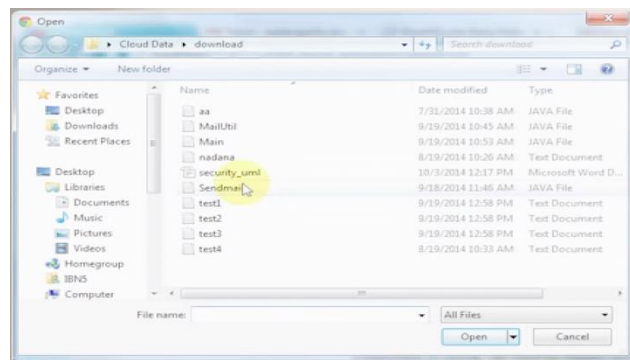


**Fig.06** Graphical File Upload

The result in Figure 6 depicts the file upload process where user is seen choosing the file to be uploaded to the cloud.
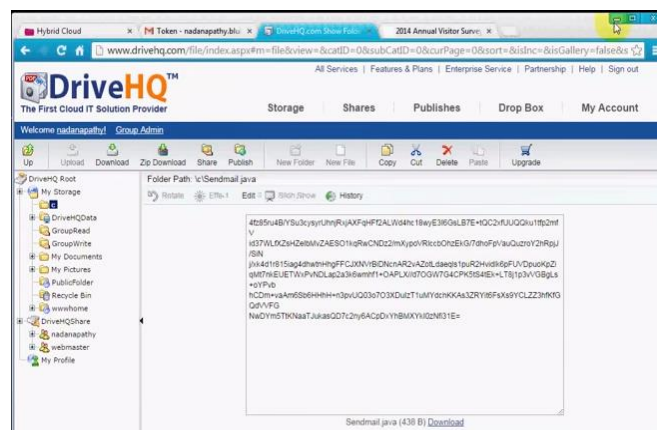


**Fig.07** File on Drive HQ

## 5. Conclusion

The prooposed research paper introduces the concept of authorized data deduplication as a means to enhance data security by incorporating differential privileges for users during the duplicate check process. The proposed approach focuses on supporting authorized duplicate check within a hybrid cloud architecture. In this setup, the private cloud server generates duplicate-check tokens for files using private keys, ensuring robust security measures.

To achieve this, several innovative deduplication constructions are presented, all of which facilitate the authorized duplicate check process in the hybrid cloud environment. The security analysis conducted on these schemes demonstrates their effectiveness in countering both insider and outsider attacks, as specified in the proposed security model. By considering potential threats from both within and outside the system, the proposed solutions provide comprehensive protection against unauthorized access and data breaches.

The research team implemented a prototype of the proposed authorized duplicate check scheme, validating its feasibility and practicality. To further assess its performance, testbed experiments were conducted on the prototype. The results show that the proposed scheme incurs minimal overhead when compared to other techniques such as convergent encryption and network transfer. This indicates that the method not only enhances data security but also maintains system efficiency, making it a promising solution for real-world applications.

Overall, this paper presents a novel and effective approach to data deduplication that prioritizes data security and addresses potential vulnerabilities in a hybrid cloud architecture. By leveraging differential privileges and utilizing private keys for generating duplicate-check tokens, the proposed schemes offer a robust defense against various types of attacks. The implementation of a prototype and the positive testbed experiment results further validate the practicality and efficiency of the approach, paving the way for its potential adoption in real-world scenarios..

## References

[1]  http://www.emc.com/leadership/digitaluniverse

[2]  J. Douceur, A. Adya, W. Bolosky, S. Dan, and M. Theimer, "Reclaiming space from duplicate files in a serverless distributed file system," in in Proceedings of the 22nd International Conference on Distributed Computing Systems, IEEE, 2002, pp. 617–624**.**

[3]  Choi [3]        C. Y. Liu ,X. J. Liu, and L. Wan, "Policy-based deduplication in secure cloud storage," in Proc. Trustworthy Comput. Serv., 2013, pp. 250–262, doi:10.1007/978-3-642-35795-4_32.

[4]  M. Bellare, S. Keelveedhi, and T. Ristenpart, "Message-locked encryption and secure deduplication," in Annual International Conference on the Theory and Applications of Cryptographic Techniques. Springer, 2013, pp. 296–312.

[5]  Y. Duan, "Distributed key generation for encrypted deduplication: Achieving the strongest privacy," in Proceedings of the 6th edition of the ACM Workshop on Cloud Computing Security. ACM, 2014, pp. 57–68.

[6]  A. Juels, and B. S. Kaliski, "PORs: Proofs of retrievability for large files," Proc. ACM Conference on Computer and Communications Security, pp. 584–597, 2007

[7]  J. Xu, E. C. Chang, and J. Zhou, "Leakage-resilient client-side deduplication of encrypted data in cloud storage," IACR ePrint Archive, 15pages, 2011.

[8]  T. Y.Wu, J. S. Pan, and C. F. Lin, Improving accessing efficiency of cloud storage using deduplication and feedback schemes,IEEESyst.J.,vol.8

[9]  C. Fan, S. Y. Huang, andW. C. Hsu, Hybrid data deduplication in cloud environment, in Proc. Int. Conf. Inf. Secur. Intell. Control,2012, pp. 174177.

[10] Y. Zhou, D. Feng, W. Xia, M. Fu, F. Huang, Y. Zhang, and C. Li, "Secdep: A user-aware efficient fine-grained secure deduplication scheme with multi-level key management," in 31st Symposium on Mass Storage Systems and Technologies (MSST), IEEE. IEEE, 2015, pp. 1–14.

[11] L. Li, X. Chen, X. Huang, S. Tang, Y. Xiang, M. M. Hassanmad, and A. Alelaiwi, "Secure distributed deduplication systems with improved reliability," IEEE Transactions on Computers, vol. 64, no. 12, pp. 3569–3579, 2015.

[12] J. Blasco, R. D. Pietro, A. Orfila, and A. Sorniotti, "A tunable proof of ownership scheme for deduplication using bloom filters," in 2014 IEEE Conference on Communications and Network Security (CNS), 2014, pp. 481–489.

[13] S. Halevi, D. Harnik, B. Pinkas, and A. Shulman-Peleg, "Proofs of ownership in remote storage systems." in ACM Conference on Computer and Communications Security, CCS 2011, Chicago, Illinois, Usa, October, 2011, pp. 491–500.

[14] A. S. R. D. Pietro, "Boosting efficiency and security in proof of ownership for deduplication," in Proceedings of the 7th ACM Symposium on Information, Computer and Communications Security, 2012, pp. 81–82.

[15] W. K. Ng, Y. Wen, and H. Zhu, "Private data deduplication protocols in cloud storage," in Proceedings of the 27th Annual ACM Symposium on Applied Computing, 2012, pp. 441–446.Elsevier (1995) 1401-1406.