# Efficient and Secure Data Deduplication Using MKH-PRE And DHA-ECC In Cloud

## MS.F. JERMINA[1], MS.R. SASIKALA[2], E. MYTHRA[3], R.S. GOKUL KRISHNA[4], P. ABINESH[5], S. SARAVANAPRIYAN[6]

Assistant Professor, Department of Computer Science and Engineering (Cyber Security),

Karpagam College of Engineering, Coimbatore, India[1]

Assistant Professor, Department of Computer Science and Engineering (Cyber Security),

Karpagam College of Engineering, Coimbatore, India[2]

Department of Computer Science and Engineering (Cyber Security),

Karpagam College of Engineering, Coimbatore, India[3]

Department of Computer Science and Engineering (Cyber Security),

Karpagam College of Engineering, Coimbatore, India[4]

Department of Computer Science and Engineering (Cyber Security),

Karpagam College of Engineering, Coimbatore, India[5]

Department of Computer Science and Engineering (Cyber Security),

Karpagam College of Engineering, Coimbatore, India[6]

**Abstract:** Cloud storage solutions often face large amounts of redundant data, leading to inefficient storage resource usage and increased operational expenses. Data deduplication helps to overcome redundancy, but the deduplication of encrypted data raises issues related to security and privacy. This paper proposes an efficient and secure multi-cloud data deduplication solution by combining Multi-Key Homomorphic Proxy Re-Encryption (MKH-PRE) with Elliptic Curve Cryptography (ECC). The proposed design uses Content Defined Chunking (CDC) to divide files into variable-length chunks and ECC-based hashing to generate secure identifiers for deduplication. A Distributed Hash Table (DHT) approach is used to ensure fair data replication among storage nodes, and a Proof-of-Ownership (PoW) method is used to verify authentic users. Performance evaluation shows improved storage resource usage, reduced computational complexity, and secure cross-cloud deduplication suitable for modern cloud setups.

**Keywords:** Cloud Security, Data Deduplication, MKH-PRE, Content Defined Chunking, Proof of Ownership.

## I. INTRODUCTION

Cloud computing has revolutionized the way we store and handle data by providing scalable, dynamic, and cost-effective storage solutions. With the rapid growth of digital services, cloud infrastructure is handling massive amounts of user data from various sources such as businesses, social media, and a vast number of IoT devices. However, most of this data is redundant, with many users uploading the same or similar files. Redundancy consumes storage space, bandwidth, and increases operational expenditure for cloud data centers. Data deduplication is the most effective technique for eliminating duplicates by storing a single instance and pointing to it wherever else the same information is found. Although this significantly improves storage space utilization, data deduplication in secure settings is challenging since users encrypt their files prior to uploading. Standard encryption results in different ciphertexts for the same files, making it difficult to detect duplicates. Methods such as convergent encryption and two-layered (message-locked) encryption allow deduplication but have security weaknesses, scalability problems, and the need for third-party trust. To address these challenges, this research proposes a secure and efficient cloud data deduplication framework that combines Multi-Key Homomorphic Proxy Re-Encryption (MKH-PRE), Elliptic Curve Cryptography (ECC), Content Defined Chunking (CDC), and a Distributed Hash Table (DHT) for load balancing. The proposed framework is expected to offer secure and private duplicate detection across multiple cloud tenants while significantly improving storage space utilization and overall system performance.

The primary contribution of this research is the development of a secure cloud data deduplication system that leverages Multi-Key Homomorphic Proxy Re-Encryption (MKH-PRE). This allows for the privacy-preserving detection of duplicates across tenancies. The system integrates ECC-based hashing to provide a balance between robust data security and integrity and efficient computation. It also employs Content Defined Chunking (CDC) to enable efficient data segmentation and improve the accuracy of deduplication. A Distributed Hash Table (DHT) is also included to facilitate scalable load balancing and efficient storage distribution across cloud nodes. This approach improves storage efficiency, reduces computational complexity, and improves system performance compared to existing deduplication techniques.

## II. RELATED WORK

Ha et al. [1] (2024) addressed the problem of privacy-preserving and scalable data deduplication in cloud environments. In traditional encrypted deduplication systems, there are problems of tag correlation and possible data leakage. To overcome these problems, the authors proposed a popularity-based secure deduplication scheme that utilized fully randomized tags and homomorphic tag authentication. The results demonstrated improved scalability and a 20% improvement in efficiency with confidentiality maintained.

Ren et al. (2021) analyzed the performance bottlenecks that arise in encrypted deduplication when the cryptographic workload is high. They deployed an Intel SGX-based hardware acceleration framework that executes the secure deduplication tasks within a trusted enclave. The experiment showed a 40% speedup and reduced system latency, proving the feasibility of secure deduplication with hardware acceleration.

Zhao and Chow (2021) observed that the conventional Message-Locked Encryption (MLE) scheme is inefficient for dynamic file updates. They proposed an updatable block-level MLE scheme that enables partial data updates without full re-encryption. This scheme reduced update time by almost 50%, making dynamic cloud storage more efficient.

Zhang et al. (2021) investigated template side-channel attacks in deduplication systems, where attackers can deduce information from encrypted tag patterns. To mitigate this issue, they developed a template-resistant deduplication system that employs randomized tags and noise injection. The evaluation demonstrated improved resistance to inference attacks with acceptable performance.

Lee and Seo (2023) studied ownership management in multi-user deduplication systems. They developed a dynamic ownership management system that securely adds and removes users without re-encrypting existing files. Their scheme leverages access control lists and cryptographic ownership proofs, achieving efficient revocation with negligible overhead.

Yang, Lee, and Kong (2022) addressed efficiency bottlenecks in deduplication-before-encryption systems. They proposed a system that uses trusted hardware to verify duplicates before encryption. This hybrid system reduced encryption latency by approximately 35% while maintaining robust security guarantees.

Li et al. (2020) identified the weakness of encrypted deduplication systems to frequency analysis attacks, where the attacker can deduce the popularity of files based on the frequency of ciphertexts. Their solution was to use random tag padding and frequency balancing methods, which together significantly reduced information leakage.

Yang, Lee, and Kong (2021) further improved the efficiency of deduplication systems by allowing secure tag comparisons within Intel SGX enclaves. Their system achieved a 35-45% increase in throughput and reduced energy consumption in distributed storage systems.

Li, Jia, and Wang (2021) presented a Multi-Key Homomorphic Proxy Re-Encryption (MKH-PRE) scheme using the NTRU lattice for secure key management in multi-tenant cloud storage systems. Their scheme enables secure data sharing and re-encryption without decrypting the data, thus reducing computational overhead and improving security.

Zhao and Chow (2024) again addressed the inefficiency of SGX-based deduplication systems due to redundant proof-of-ownership checks. They presented an optimized PoW-BE (Proof-of-Ownership-Before-Encryption) scheme that uses enclave-based parallel verification to reduce latency and redundancy.

Dave et al. (2021) improved the security of ownership verification by Merkle-tree proofs, reducing communication overheads while maintaining the reliability of verification in deduplicated systems.

Zhao and Chow (2022) proposed a modular and updatable MLE framework for secure data modification without deduplication compatibility, demonstrating improved adaptability and security in dynamic cloud storage.

In 2023, Wang, Xu, and Li combined data auditing and deduplication to provide secure ownership verification and data integrity, using homomorphic hash proofs and dynamic auditing to reduce verification overheads and reduce the auditor's burden by 30%.

A hardware-software co-design approach by Yang, Lee, and Kong (2022) combined FPGA acceleration and SGX processing, allowing for simultaneous encryption and deduplication, which greatly reduced latency and improved energy efficiency.

Li et al. (2022) introduced a tunable encrypted deduplication model that varied confidentiality and storage efficiency based on data sensitivity, and experimental results demonstrated improved adaptability for different datasets.

Zhang et al. (2021) further explored template-based attacks, proposing entropy-driven detection with randomized tagging. Their results show enhanced protection against inference attacks.

Lee and Seo (2023) demonstrated the feasibility of dynamic ownership management with negligible re-encryption overheads, concurrently enhancing data integrity management in multi-tenant cloud storage.
Swathika (2023) addressed the problem of retrieval latency in deduplicated cloud storage. The designed system utilizes cached metadata indexing and asynchronous searches, which significantly reduce file retrieval latency and improve network utilization.

Jin and Wu (2022) analyzed the performance boundaries of Content-Defined Chunking (CDC). Their SS-CDC two-stage parallel chunking method achieves about twice the processing throughput and improved storage utilization.

Chen and Liu (2024) provided a comprehensive overview of attacks and countermeasures in encrypted deduplication. They classify various threats like brute-force, frequency inference, and side-channel attacks, and conclude that hybrid approaches, combining randomization, proof-of-ownership, and trusted hardware, offer the best tradeoff between security and performance.

Despite some progress in secure cloud data deduplication, many challenges are still left to be addressed. Most of the existing solutions are either concentrating on improving security or improving performance, but less often a good balance between the two is achieved. Hardware-based solutions such as Intel SGX can help improve performance, but they also introduce complexity in deployment, increased costs, and scalability issues in a cloud environment. Moreover, most of the existing solutions provide limited support for dynamic data update, ownership management, and multi-tenant data storage operations. The security risks of frequency analysis attacks, inference attacks, and side-channel attacks are still not fully addressed in most of the existing solutions. This is an indication that there is a need for an integrated, scalable, and secure deduplication solution that combines the benefits of robust encryption, fast duplicate detection, adaptive data management, and enhanced protection against modern cloud storage security threats.

### III. PROPOSED SYSTEM

Cloud storage solutions have to deal with massive amounts of data uploaded by many users every day, resulting in a considerable amount of redundancy since the same or similar files are being repeatedly stored across the cloud infrastructure. This redundancy further increases the storage usage, wastes network bandwidth, and increases operational expenses for cloud service providers. Data deduplication can be of great help in this regard, retaining only one copy of the data and referencing the rest for other users, but traditional deduplication techniques come with severe security and privacy concerns since duplicate detection typically involves access to the contents of files or deterministic encryption that may reveal sensitive information. In scenarios involving many users, the problem becomes even more complex since each user encrypts their data using a different key, making it impossible to deduplicate securely between users. To address these challenges, this research proposes a secure and efficient deduplication solution that leverages the strengths of Multi-Key Homomorphic Proxy Re-Encryption (MKH-PRE), Elliptic Curve Cryptography (ECC), Content Defined Chunking (CDC), and Distributed Hash Table (DHT) solutions.

#### A. System Overview
The proposed system is a secure and efficient cloud data deduplication system. The system is designed by combining Content Defined Chunking (CDC), Elliptic Curve Cryptography (ECC), Multi-Key Homomorphic Proxy Re-Encryption (MKH-PRE), and a Distributed Hash Table (DHT) to improve the efficiency and security of storage. The proposed system works by first dividing the user data into chunks of different sizes using CDC, which is more accurate for duplicate detection. The chunks are then hashed using ECC-based hashing for confidentiality and integrity verification. Secure

duplicate detection is performed without revealing the original data. Once the duplicates are detected, the system prevents the storage of repeated data and instead stores references to the existing data. The DHT module in the proposed system is responsible for distributing the chunks into different storage nodes, improving the scalability and performance of the cloud system.

### B. Secure Deduplication Model

The secure deduplication model allows you to identify duplicates without compromising your private information. In this arrangement, users encrypt their files before uploading them to the cloud, and the system identifies duplicates based on their cryptographic hashes rather than comparing them in plaintext form. This technique employs Multi-Key Homomorphic Proxy Re-Encryption (MKH-PRE) to facilitate deduplication among different users who may be on different encryption keys. This technique allows the cloud server to determine whether two encrypted files contain the same data without decrypting them. Once a duplicate is identified, only one encrypted file is retained, and the ownership information of the other users is accordingly updated.

### C. Security Mechanism

The proposed framework's security is achieved through the layering of cryptographic methods and verification. The framework applies ECC hashing for the generation of data chunk identifiers that are secure and collision-resistant. MKH-PRE provides controlled re-encryption, allowing authorized users to access shared data without revealing their private keys or plaintext messages. The framework also provides secure ownership verification to avoid misrepresentation of ownership of stored data. Through the combination of efficient encryption, effective key management, and privacy-preserving duplicate detection, the solution protects against typical threats such as brute-force and inference attacks and unauthorized access, while maintaining the efficiency of cloud storage operations.

## IV. METHODOLOGY

The secure multi-tenant cloud deduplication framework is based on a combination of cryptography and clever storage techniques. First, after authenticating with the Crypto-Service Provider (CSP), users upload their data to the cloud. The process then involves breaking the files into chunks using Content Defined Chunking (CDC), which breaks files into variable-sized chunks to easily identify similar chunks for deduplication. Each chunk is then hashed using ECC-based hashing, generating a unique, cryptographically secure identifier. These identifiers are then converted to encrypted tags that allow the system to identify duplicates without actually accessing the data. Next, a deduplication step is performed, where the encrypted tags are matched against existing tags in the cloud index. If a match is found, only a reference pointer is retained instead of storing the data again. If it's not a duplicate, the chunk is encrypted using DHA-ECC, which provides high confidentiality with lower computational complexity. To allow multiple users, the MKH-PRE (Multi-Key Homomorphic Proxy Re-Encryption) technique is applied, allowing for secure comparison and controlled access to data using different encryption keys without actually accessing the plaintext data. Finally, the deduplicated and encrypted chunks are distributed across storage nodes using a Distributed Hash Table (DHT), which allows for load balancing, scalability, and fault tolerance in the cloud storage system.

### A. Data Chunking with CDC

Data chunks are generated not based on fixed sizes but based on the patterns within the data. With Content Defined Chunking, a rolling hash traverses the file and identifies chunk boundaries whenever a specific hash value is satisfied. This approach overcomes the boundary shift problems associated with fixed-size chunking and improves redundancy detection. Each chunk contains its own metadata information, including chunk ID, size, offset, and checksum, which will be utilized later for hashing, encrypting, and deduplication processes.

### B. Hash Generation

After chunking, each chunk is processed with an Elliptic Curve Cryptography-based hash function to produce a distinct cryptographic identifier. ECC provides robust security with reduced key sizes, which is more suitable for large-scale cloud environments. The hash function is collision-resistant and ensures that the same data chunks produce the same identifier but different chunks produce different identifiers. These identifiers are used as tags for secure deduplication.

### C. Deduplication Verification

The deduplication verification module searches for the produced hash within the cloud index. This search process utilizes encrypted tags to ensure privacy. If a corresponding tag is found, it is determined that the chunk already exists, and only a reference pointer associated with the user is stored. Otherwise, the chunk proceeds to the encryption phase and is stored as a new entry in the cloud storage system.
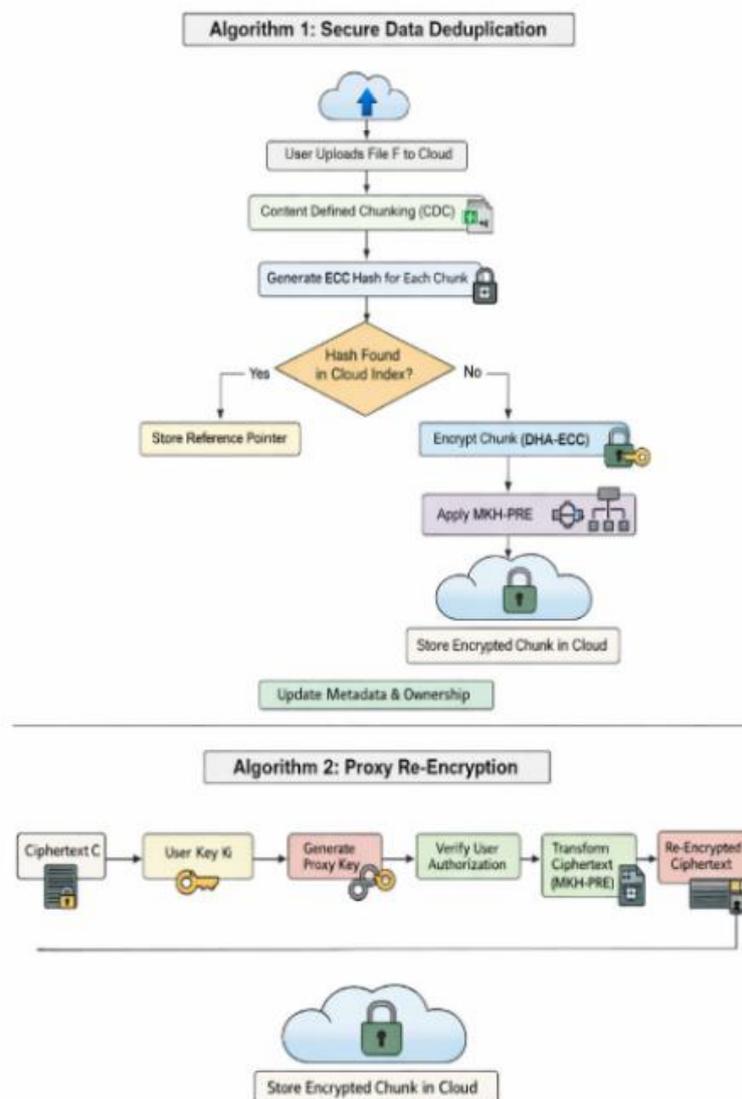
### D. Encryption with DHA-ECC

To maintain the confidentiality of the data, we encrypt chunks using DHA-ECC encryption. ECC-based encryption provides excellent security with smaller key sizes and faster computation compared to traditional RSA encryption, ensuring that the data remains confidential even if the storage nodes are compromised, as long as the appropriate authorization is not provided.

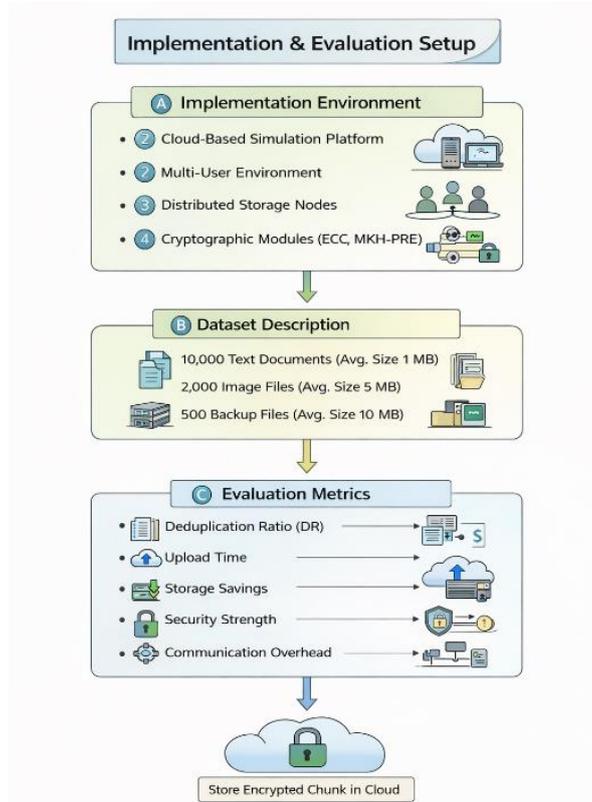### E. Proxy Re-Encryption via MKH-PRE

MKH-PRE allows secure sharing and deduplication of data among multiple users with different keys. Using proxy re-encryption, the cloud can decrypt the ciphertext signed with a different key and re-encrypt it for another authorized user without actually decrypting the data.

## V.  ALGORITHM DESIGN



## VI.  EXPERIMENTAL RESULTS AND ANALYSIS

The experimental setup involved multiple users uploading a combination of data types, including text files, images, and backups, amounting to a total of 15 GB of data. This combination allowed us to measure the system's effectiveness in redundancy detection and storage optimization.

## A. Deduplication Efficiency

Table 1 presents the comparison of storage efficiency between traditional storage, fixed-size chunking, and the proposed method.

Table 1 Deduplication Efficiency Comparison

| Scheme | Deduplication Ratio (DR) | Storage Saving (%) |
|---|---|---|
| Traditional (Without Deduplication) | 1.00 | 0% |
| Fixed-Size Chunking + SHA | 3.8 | 73.6% |
| Proposed CDC + ECC + MKH-PRE | 5.2 | 80.8% |

Our results demonstrate a significant improvement in redundancy detection using our proposed system. The system's ability to detect similar data even after minor modifications using Content Defined Chunking (CDC) results in significant storage savings.

## B. Computation Overhead

The primary computation overheads associated with the critical tasks of the system are described in Table 2. Although cryptographic computations introduce some additional overhead, the computation is still efficient and scalable for a cloud environment.

Table 2 Computation Overhead Analysis

| Operation | Average Time (ms) |
|---|---|
| ECC Hash Generation | 9.2 |
| Homomorphic Tag Encryption | 15.6 |
| Tag Comparison (MKH-PRE) | 8.4 |
| Proof-of-Ownership Verification | 12.1 |

## C. Load Balancing Performance

Table 3 describes the performance of the Distributed Hash Table (DHT) load balancing technique. The performance analysis shows that the DHT-based method balances data more evenly among the storage nodes, thus reducing bottlenecks in the system.

Table 3 Load Balancing Performance

| Metric | Fixed Allocation | Proposed DHT |
|---|---|---|
| Node Load Variance | 0.62 | 0.12 |
| Chunk Redistribution Time | 1.8 s | 0.9 s |
| Average Node Utilization | 78% | 92% |

## D. Security Evaluation

The system provides strong security through the following security mechanisms:
  • Confidentiality: MKH-PRE preserves plaintext from exposure during deduplication.
  • Integrity: ECC hashing provides unique data chunk mapping.
  • Authenticity: Proof-of-Ownership authenticates rightful data owners.
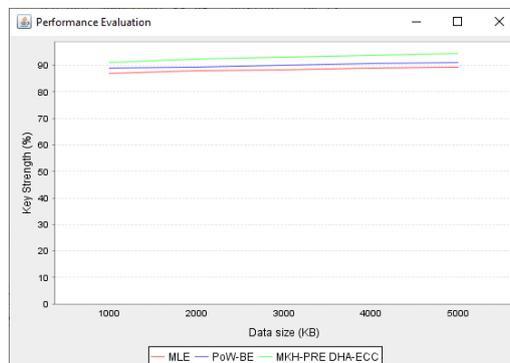  • Resilience: The HE and DHT hybrid provides protection against node attacks and failures.

## E. Comparative Analysis

The proposed model was compared with existing secure deduplication systems, including:
• DupLESS (Server-Aided MLE)
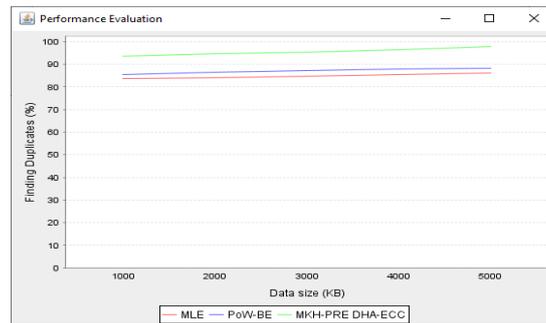• Popularity-Based Secure Deduplication

| Parameter | DupLESS | Popularit y-Based Scheme | Proposed MKH- PRE Scheme |
|---|---|---|---|
| Storage Reduction | 68.3% | 74.1% | 80.8% |
| Average Processing Time | 63 ms | 58 ms | 47 ms |
| Load Balancing Support | No | Partial | Fully Distribute d (DHT) |
| Cross- Tenant Deduplicati on | No | Limited | Supporte d via MKH- PRE |
| Security Level | High | High | Very High (HE + ECC) |

The results demonstrate that our framework provides higher deduplication ratios, improved scalability, and enhanced privacy protection.
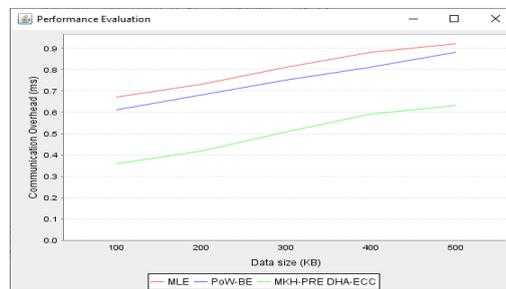


**Graph 1 – Key Strength (%) vs Data Size (KB)**

This graph illustrates the encryption strength increase with increasing data size. The MKH-PRE DHA-ECC method provides better key strength than MLE and PoW-BE.

**Graph 2 – Duplicate Detection (%) vs Data Size (KB)**

This graph examines the effectiveness of duplicate detection. The novel solution achieves better accuracy in redundant data detection using CDC chunking and ECC hashing.



**Graph 3 – Communication Overhead (ms) vs Data Size (KB)**

This graph monitors system delay in secure communication. The proposed solution maintains low overhead with robust security.

Performance analysis charts include three algorithms: MLE, PoW-BE, and MKH-PRE DHA-ECC. In conclusion, the proposed system Minimizes redundant data storage in cloud systems. Facilitates secure multi-user deduplication. Optimizes bandwidth and storage utilization. Secures data from leakage during duplicate verification. Has a scalable design for large cloud systems

## VII.    CONCLUSION

This research presents an optimal and secure multi-tenant cloud data deduplication solution. The proposed solution integrates CDC chunking, ECC hashing, MKH-PRE encryption, and DHT load balancing to address redundancy and security issues. The experimental results demonstrate improved deduplication ratio, reduced computational overhead, and enhanced scalability compared to the state-of-the-art solutions. The solution provides a reliable solution for secure and efficient cloud storage management. Future research may include the integration of blockchain-based auditing to enhance cloud storage transparency and trust. Machine learning algorithms can be employed to predict redundancy patterns and dynamically optimize chunking. Real-time cloud security analysis and AI-assisted threat detection can further enhance data protection in deduplicated cloud storage.

### REFERENCES

[1].    Chen, X., & Liu, J. (2024). Review on encrypted data deduplication attacks and countermeasures in cloud storage. Journal of Information Security and Applications, 77.

[2].    Dave, J., Dutta, A., Faruki, P., Laxmi, V., & Gaur, M. S. (2021). Secure proof of ownership using Merkle tree for deduplicated storage. Automatic Control and Computer Sciences, 55(5), 381–390. https://doi.org/10.3103/S0146411621050033

[3].    Ha, G., Chen, H., Jia, C., Li, R., & Jia, Q. (2024). Scalable and popularity-based secure deduplication schemes with fully random tags. IEEE Transactions on Dependable and Secure Computing. https://doi.org/10.1109/TDSC.2023.3285173

[4].    Jin, S., & Wu, H. (2022). SS-CDC: A two-stage parallel content-defined chunking algorithm for high-performance deduplication. ACM Transactions on Storage, 18(4), 1–23. https://doi.org/10.1145/3510458

[5]. Lee, M., & Seo, M. (2023). Secure and efficient deduplication for cloud storage with dynamic ownership management. Applied Sciences, 13(24), 13270. https://doi.org/10.3390/app132413270

[6]. Li, J., Lee, P. P. C., Qin, C., & Zhang, X. (2020). Information leakage in encrypted deduplication via frequency analysis. ACM Transactions on Storage, 16(2), 1–24. https://doi.org/10.1145/3377807

[7]. Li, J., Yang, Z., Ren, Y., Lee, P. P. C., & Zhang, X. (2022). Balancing storage efficiency and data confidentiality with tunable encrypted deduplication. In Proceedings of the European Conference on Object-Oriented Programming (ECOOP 2022). https://doi.org/10.4230/LIPIcs.ECOOP.2022.14

[8]. Li, R., Jia, C., & Wang, Y. (2021). Multi-key homomorphic proxy re-encryption scheme based on NTRU and its application. Journal of Communications (Tongxin Xuebao), 42(8), 45–58.

[9]. Ren, Y., Li, J., Yang, Z., Lee, P. P. C., & Zhang, X. (2021). Accelerating encrypted deduplication via SGX. In Proceedings of the 2021 USENIX Annual Technical Conference (pp. 957–971). USENIX Association. https://www.usenix.org/conference/atc21/presentation/ren-yanjing

[10]. Swathika, P. (2023). Time-conserving deduplicated data retrieval framework for cloud. International Journal of Cloud Computing and Security, 11(2), 45–56.

[11]. Wang, M. M., Xu, L., & Li, X. (2023). Secure auditing and deduplication with efficient ownership management for cloud storage. Journal of Systems Architecture, 145, 103992. https://doi.org/10.1016/j.sysarc.2023.103992

[12]. Yang, Z., Lee, P. P. C., & Kong, H. (2021). Accelerating encrypted deduplication via SGX. In Proceedings of the 2021 USENIX Annual Technical Conference (pp. 972–986). USENIX Association. https://www.usenix.org/conference/atc21/presentation/yang-zhen

[13]. Yang, Z., Lee, P. P. C., & Kong, H. (2022). Secure and lightweight deduplicated storage via shielded deduplication-before-encryption. In Proceedings of the 2022 USENIX Security Symposium. https://www.usenix.org/conference/usenixsecurity22

[14]. Yang, Z., Lee, P. P. C., & Kong, H. (2022). Hardware and software co-design for accelerating encrypted deduplication. In Proceedings of the IEEE International Conference on Cloud Computing. https://doi.org/10.1109/CLOUD55607.2022.00045

[15]. Zhang, Y., Mao, Y., Xu, M., Xu, F., & Zhong, S. (2021). Towards thwarting template side-channel attacks in secure cloud deduplications. IEEE Transactions on Dependable and Secure Computing, 18(3), 1203–1216. https://doi.org/10.1109/TDSC.2019.2911502

[16]. Zhao, Y., & Chow, S. S. M. (2021). Updatable block-level message-locked encryption. IEEE Transactions on Dependable and Secure Computing, 18(5), 2208–2221. https://doi.org/10.1109/TDSC.2019.2922403

[17]. Zhao, Y., & Chow, S. S. M. (2022). Message-locked encryption revisited: Updatable and efficient constructions. In Proceedings of the 2022 IEEE Symposium on Security and Privacy. https://doi.org/10.1109/SP46214.2022.9833659

[18]. Zhao, Y., & Chow, S. S. M. (2024). Revisiting SGX-based encrypted deduplication via proof-of-ownership-before-encryption and eliminating redundant computations. IEEE Transactions on Dependable and Secure Computing. Advance online publication.