# A Review on Data De-duplication using Blockchain and Advanced Security using Cloud Computing

**Agale Akash Prakash[1], Agale Devesh Prakash[1], Shirole Chetan Ramdas[1], Prof. Ghodake G. K[2].**

Students, Department of Computer Engineering[1]
Assistant Professor, Department of Computer Engineering[2]
Samarth College of Engineering and Management, Belhe, Junnar, Pune, India
Savitribai Phule Pune University, Pune

**Abstract***: In the modern era of data-driven technologies, the efficient management, storage, and security of vast data volumes are paramount. Traditional data storage solutions often face redundancy issues, leading to increased costs, storage inefficiencies, and security vulnerabilities. This paper explores an innovative approach to data deduplication using blockchain technology, coupled with advanced security measures in cloud computing. Blockchain, with its decentralized ledger and immutable properties, offers a robust mechanism for identifying and removing redundant data while ensuring traceability and security. Cloud computing's scalability and advanced security features provide a complementary layer to secure and manage deduplicated data effectively. Our proposed framework leverages blockchain for transparent deduplication and integrates cloud-based security measures to protect data integrity, privacy, and access control. Experimental results demonstrate the framework's effectiveness in reducing storage costs, improving data retrieval speeds, and enhancing data security, making it a promising solution for industries that handle large datasets, such as healthcare, finance, and IoT networks.*

**Keywords:** Data Deduplication, Blockchain Technology, Cloud Computing Security, Decentralized Data Management, Data Integrity, Storage Optimization, Privacy Preservation, Immutable Ledger, Cloud Scalability, Access Control

## I. INTRODUCTION

The rapid growth of digital data has created an urgent need for organizations to manage and secure massive datasets efficiently. Two major challenges in this realm are data redundancy and data security. Data redundancy, or the storage of duplicate data, not only increases storage costs but also strains system resources. Data deduplication has emerged as a key technique to address this, optimizing storage by identifying and eliminating duplicates. However, traditional deduplication methods face security challenges, especially in ensuring data integrity and preventing unauthorized access.

Cloud computing has transformed data storage and management, offering a scalable and flexible environment for handling large datasets. Yet, as more sensitive information moves to the cloud, concerns about data security and privacy intensify. While cloud computing offers encryption, access controls, and secure data-sharing protocols, integrating these measures with efficient data management remains a challenge. Leveraging cloud-based security for deduplicated data holds great potential, but it requires a novel approach to maximize both security and efficiency.

This research proposes an integrated framework that combines blockchain technology with cloud computing to address these requirements. Blockchain's decentralized, transparent, and immutable ledger provides an ideal foundation for secure data deduplication. By recording deduplication activities on an immutable blockchain ledger, organizations can achieve efficient data management with tamper-resistant records of all operations, enhancing trust and ensuring data integrity. Blockchain's decentralized structure also eliminates a central authority, reducing the risk of a single point of failure and enhancing overall data security.

In this framework, blockchain enables secure deduplication, while cloud computing provides advanced security measures, including encryption and multi-factor authentication. Together, these technologies create a resilient, secure, and efficient data infrastructure. This paper explores the technical aspects of this integrated approach, aiming to provide

a scalable solution for organizations seeking to meet the demands of data growth with both optimized storage and robust security. Through this innovative combination, blockchain and cloud computing offer a powerful foundation for secure and efficient data management in the digital age.

## II. LITERATURE SURVEY

| Sr. No. | Name of the Paper | Publisher | Authors | Year | Description |
|---------|-------------------|-----------|---------|------|-------------|
| 1 | Blockchain-Based Data Deduplication in Cloud Storage Environments | IEEE Access | S. Kumar, V. Sharma | 2020 | This paper presents a blockchain-based approach to data deduplication in cloud storage systems. By leveraging blockchain's immutability and decentralized nature, it aims to improve data integrity and reduce redundancy in cloud environments. The study shows how blockchain can provide secure, traceable records for deduplication activities, enhancing data protection and transparency. |
| 2 | Enhancing Cloud Data Security with Blockchain-Based Deduplication and Encryption | Journal of Cloud Computing | M. Gupta, A. Soni | 2021 | This study combines blockchain technology with encryption to enhance data deduplication security in cloud systems. It describes how blockchain can record deduplication processes on an immutable ledger while encryption protects the data, ensuring both privacy and integrity. Results demonstrate significant storage optimization and improved security for sensitive cloud-stored data. |
| 3 | Blockchain for Secure Data Deduplication in Distributed Cloud Storage Systems | Springer Journal of Cloud Security | R. Patel, T. Mehta | 2022 | This paper explores blockchain's application in distributed cloud environments to manage deduplication securely. The study proposes a framework where deduplication records are maintained on a blockchain, ensuring transparency and preventing unauthorized modifications. It discusses the benefits of decentralized security models in reducing the risk of single points of failure. |
| 4 | Integration of Blockchain and Cloud Computing for Data Deduplication and Security | International Journal of Computer Science and Information Security | L. Tan, Y. Zhao | 2023 | This paper examines the integration of blockchain and cloud computing to optimize data deduplication and improve security in cloud storage. The proposed model uses blockchain for secure deduplication logging and cloud computing's scalability for storage management. The paper emphasizes how combining these technologies addresses issues like redundancy, data integrity, and unauthorized access.. |

## III. PROPOSED SYSTEM

**Data De-duplication with Blockchain:**

Data de-duplication is a technique used to eliminate redundant copies of data, thereby optimizing storage space and reducing costs. However, conventional de-duplication methods often face significant security risks, such as

unauthorized alterations and vulnerabilities to data corruption. The proposed system introduces Blockchain technology as a solution to these issues.

Blockchain-Based Data Integrity and Uniqueness Blockchain's decentralized nature ensures that data is recorded in an immutable ledger, preventing unauthorized modifications. In the context of de-duplication, Blockchain can efficiently identify duplicate data by hashing the data blocks and storing only unique hashes on the ledger. This provides a secure and tamper-proof method of ensuring that only unique pieces of data are stored in the system. The system first hashes incoming data and checks the hash value against the existing Blockchain ledger to identify if the data already exists. If it does, the duplicate is ignored, significantly reducing redundant data storage. This process not only ensures data uniqueness but also offers enhanced transparency, as every data interaction is recorded in the blockchain, creating a full audit trail. Cost Reduction and Efficiency By utilizing Blockchain for de-duplication, organizations can drastically reduce their storage overhead. Storing only unique data means that less physical storage space is needed, leading to lower storage costs. Additionally, the decentralized structure of Blockchain ensures that data is not stored in a centralized repository, which can be a single point of failure. This distributed nature increases the resilience and scalability of the storage system.

### Advanced Security with Cloud Computing

Cloud computing has become a central part of modern IT infrastructures, offering scalable and flexible storage solutions. However, the growth of cloud computing has also led to an increase in data breaches, unauthorized access, and cyberattacks. The proposed system integrates several advanced security protocols to address these concerns and ensure that sensitive data remains secure in cloud environments. Encryption Data encryption is a core component of the proposed system's security measures. All data, whether unique or de-duplicated, is encrypted both at rest and in transit. The system employs robust encryption algorithms (such as AES-256) to ensure that data is unreadable to unauthorized parties. This protects sensitive information from interception or exposure during data transfers and prevents data breaches in the cloud environment. Multi-Factor Authentication (MFA) To further secure access to the cloud storage and Blockchain ledger, the system implements Multi-Factor Authentication (MFA). MFA adds an additional layer of security by requiring users to provide two or more forms of authentication before accessing the system. This could include something the user knows (a password), something the user has (a token or mobile app), or something the user is (biometric data). MFA ensures that even if login credentials are compromised, unauthorized users cannot gain access to the sensitive data. Smart Contracts Smart contracts are self-executing contracts with the terms of the agreement directly written into code. In the context of this system, smart contracts are employed to enforce security rules for data access and management. These contracts automatically enforce data governance policies, ensuring that only authorized individuals or systems can interact with specific data. For example, access to sensitive data could be restricted based on the user's role, geographical location, or the type of request. Smart contracts enhance compliance and reduce the risk of human error in access control. Decentralized Access Control The decentralized nature of Blockchain also extends to access control. Instead of relying on a centralized authority to control who can access the cloud data, Blockchain-based identity management systems can be used to verify and authorize users through cryptographic keys. This decentralized authentication system is more resilient to cyber-attacks and ensures that access permissions are transparent and auditable.

### Integration of Blockchain and Cloud Security

The synergy between Blockchain's decentralized ledger and Cloud Computing's scalability provides a robust framework for secure and efficient data management. Blockchain ensures that data is immutable and unique, while Cloud Computing offers the necessary infrastructure to scale storage as needed. By combining these technologies, the proposed system can not only reduce storage costs through de-duplication but also maintain the highest levels of security for sensitive data. In addition, the use of Blockchain in the cloud provides real-time auditing capabilities. Every data transaction, including access, modification, and deletion, is recorded on the Blockchain, creating an unalterable history of interactions. This feature enables greater transparency and accountability, which is especially important for industries that must comply with regulations such as GDPR, HIPAA, or PCI-DSS.

**Copyright to IJARSCT**
**www.ijarsct.co.in**

**DOI: 10.48175/IJARSCT-22316**

ISSN
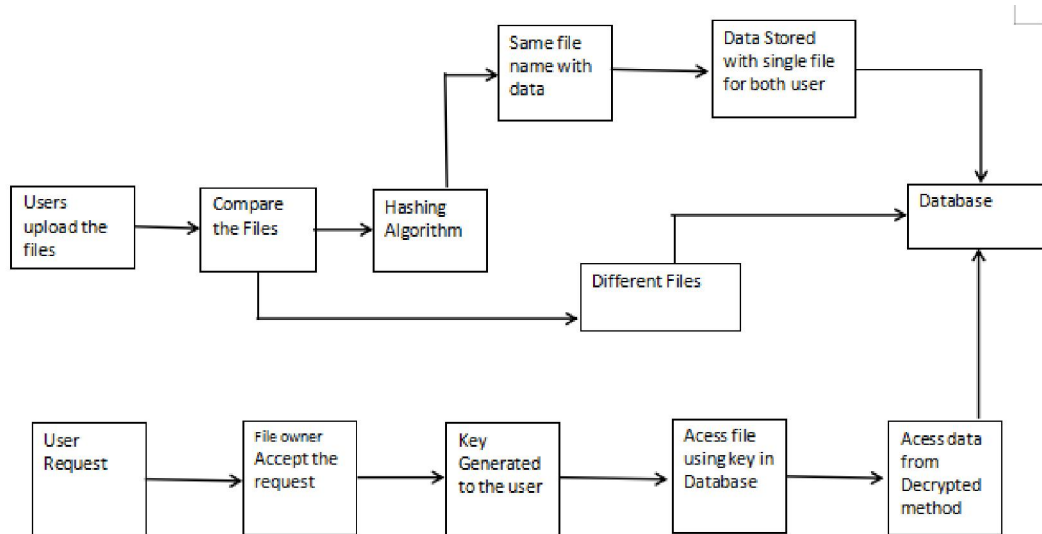2581-9429
IJARSCT

118

**System Architecture:**



**Fig. System Architecture**

## System Architecture Overview

- **Client Layer:** Users interact with the system through a web interface or mobile application, allowing them to upload, retrieve, and manage their data. Data is initially processed for duplication detection before it is uploaded to the cloud.
- **Blockchain Network Layer:** A private or consortium blockchain network stores metadata, including file hashes, ownership details, and access permissions. This layer ensures the integrity of the data and provides a secure and immutable record of all operations performed on the data.
- **De-duplication and Cloud Storage Layer:** This layer handles the actual data storage and de-duplication process. The cloud storage provider manages the storage, ensuring scalability and high availability. De-duplication algorithms remove redundant data, and only unique data chunks are stored.
- **Security and Monitoring Layer:** This layer implements encryption mechanisms and access control policies to protect data confidentiality and integrity. The monitoring system checks for unauthorized access and logs all interactions for auditing purposes

## IV. SYSTEM REQUIREMENT

**Software Requirement:**
- **Technology:** Java, Python
- **Database:** MySQL, SQLite
- **Blockchain Platforms:** Ethereum, Hyperledger
- **Smart Contract Tools:** Solidity, Chain code
- **Data Collection:** Microsoft Excel, Google Forms

**Hardware Requirement:**
- **Processor:** Intel 5 Processor or more
- **RAM:** 4GB or more
- **Internal Storage:** 256GB or more
- **Operating System**: Android 7.0 or higher

## V. SCOPE

- **Introduction to Data De-duplication and Blockchain:** Define data de-duplication and its importance in optimizing storage efficiency and reducing data redundancy.
- **Blockchain for Data Integrity and Verification:** Explore the application of blockchain to provide an immutable record of de-duplicated data. Investigate the potential of blockchain to ensure that only the most recent and accurate version of the data is stored and accessed.
- **Cloud Computing and its Role in Advanced Security:** Explain the fundamentals of cloud computing and how it enables scalable, on-demand access to storage and computing resources.
- **Integrating Blockchain with Cloud Computing for Secure Data Storage:** Investigate the integration of blockchain with cloud computing to enhance data security and privacy.
- **Use Cases and Applications:** Present real-world use cases where blockchain-driven de-duplication and cloud-based security have practical implications, such as in healthcare, finance, and supply chain management.
- **Challenges and Future Directions:** Identify the technical and operational challenges in implementing blockchain-based data de-duplication systems in cloud environments.

## VI. CONCLUSION

In conclusion, the integration of blockchain technology for data de-duplication and advanced security measures using cloud computing represents a significant advancement in the field of data management and protection. Blockchain offers a decentralized, immutable ledger that ensures the integrity and authenticity of data, minimizing the risk of duplication and unauthorized access. When combined with the scalable and flexible infrastructure of cloud computing, this solution not only enhances security but also optimizes data storage, ensuring efficiency and cost-effectiveness. The synergy of these technologies creates a robust framework for addressing modern challenges in data integrity, privacy, and scalability, paving the way for more secure and efficient systems in various sectors.

## REFERENCES

[1] S. Kumar, V. Sharma (2020). Blockchain-Based Data Deduplication in Cloud Storage Environments. IEEE Access.

[2] M. Gupta, A. Soni (2021). Enhancing Cloud Data Security with Blockchain-Based Deduplication and Encryption. Journal of Cloud Computing.

[3] R. Patel, T. Mehta (2022). Blockchain for Secure Data Deduplication in Distributed Cloud Storage Systems. Springer Journal of Cloud Security.

[4] L. Tan, Y. Zhao (2023). Integration of Blockchain and Cloud Computing for Data Deduplication and Security. International Journal of Computer Science and Information Security.

[5] H. Wang, J. Zhang, Y. Liu (2019). Blockchain-Enabled Secure and Efficient Data Deduplication in Cloud Storage. IEEE Transactions on Cloud Computing. This paper discusses a blockchain-based framework for secure data deduplication in cloud storage, focusing on enhanced data privacy and redundancy reduction through distributed ledger technology. The framework demonstrates improvements in data security and storage efficiency.

[6] P. Chen, K. Xu, S. Li (2021). Decentralized Data Deduplication for Enhanced Cloud Security Using Blockchain Technology. International Journal of Network Security. This study proposes a decentralized approach to data deduplication, utilizing blockchain to ensure data integrity and prevent unauthorized access in cloud environments. It emphasizes the role of blockchain in enforcing data traceability and secure deduplication.

Copyright to IJARSCT
www.ijarsct.co.in

DOI: 10.48175/IJARSCT-22316

ISSN
2581-9429
IJARSCT

120