# Integrating Blockchain with Big Data for Secure Data Sharing: A Comprehensive Methodology

**Akash Hooda[1*], Arju Hooda[2] and Disha Yadav[3]**

[1]*Netaji Subhas Institute of Technology, Dwarka, New-Delhi*

[2]*Delhi Technological University, Rohini, New-Delhi*

[3]*Indira Gandhi Delhi Technical University for Women, Kashmere Gate, New-Delhi*

***Corresponding Author**
Akash Hooda, Netaji Subhas Institute of Technology, Dwarka, New-Delhi.

## Abstract

*In today's data-driven world, the volume of data managed by organizations is growing rapidly, presenting significant challenges in ensuring data security, integrity, and scalability. While blockchain technology offers robust security features, such as immutability and decentralization, it struggles with scalability issues, particularly in high-throughput environments. Conversely, big data frameworks like Hadoop and Spark excel in handling large datasets efficiently but often lack strong security mechanisms.*

*This research proposes a hybrid architecture that integrates the security of blockchain with the scalability of big data frameworks, creating a system capable of securely managing vast amounts of data in real-time. The architecture includes advanced encryption methods, off-chain data management, and seamless integration with existing big data tools, making it suitable for industries such as healthcare, finance, and IoT. Through a comprehensive methodology involving literature review, requirement analysis, architectural design, and performance evaluation, the study demonstrates that this hybrid approach significantly enhances both security and scalability, offering a future-ready solution for secure data sharing across various sectors.*

## 1. Introduction

As organizations continue to generate and process massive amounts of data, they face the dual challenge of ensuring data security while maintaining scalability. Big data frame- works like Hadoop and Spark are designed to manage large datasets efficiently, but they often fall short in providing the robust security needed to protect sensitive information. On the other hand, blockchain technology, with its decentralized and immutable nature, offers strong security but struggles with scalability, especially when dealing with high transaction volumes and large datasets.

This research focuses on developing a hybrid architecture that integrates the security strengths of blockchain with the scalability of big data frameworks. The goal is to create a system that can securely manage and share large volumes of data in real-time, meeting the needs of industries where data integrity, privacy, and compliance are critical, such as health- care, finance, and supply chain management. By designing a system that balances these two technologies, we aim to provide a solution that not only addresses current challenges but is also adaptable to future demands.

## 2. Problem Statement

The integration of blockchain technology with big data plat- forms offers significant potential for secure and scalable data sharing across organizations, particularly in high-throughput environments. However, several challenges hinder the full realization of this potential:

### 2.1 Encryption and Data Security
### 2.1.1 Complexity

Implementing encryption within blockchain-big data systems is inherently complex, especially when dealing with large datasets that require fast, secure access. The need to balance robust encryption with efficient decryption processes for legitimate users presents significant technical challenges.

### 2.1.2 Privacy

Maintaining data confidentiality in blockchain's transparent

environment is difficult, particularly in scenarios where data is shared among multiple parties. The inherent transparency of blockchain can conflict with the need to protect sensitive information, necessitating the development of advanced privacy-preserving techniques, such as homomorphic encryption and zero-knowledge proofs.

## 2.2 Key Management
### 2.2.1 Scalability
Managing cryptographic keys effectively in large, distributed networks is a major challenge. Traditional key management solutions often struggle to scale in such environments, leading to potential vulnerabilities and inefficiencies.

### 2.2.2 Decentralization
Distributing cryptographic keys securely without relying on a central authority remains an unresolved issue in decentralized systems. This challenge is particularly pronounced in blockchain networks, where ensuring that all participants have secure access to the necessary keys without introducing single points of failure is crucial.

## 2.3 Scalability
### 2.3.1 Storage Overhead
Blockchain technology introduces significant storage overhead due to the need to maintain a continuously growing ledger that records every transaction. In the context of big data, where datasets are frequently updated, this can lead to substantial storage requirements that slow down data processing and retrieval.

### 2.3.2 Throughput and Latency
Blockchain systems, especially those utilizing consensus mechanisms like Proof of Work (Po W) or Proof of Stake (Po S), have limited trans- action throughput and can introduce significant latency. These limitations are particularly problematic in high-throughput environments, such as real-time data analytics or financial services, where large volumes of transactions must be processed quickly.

## 2.4 Privacy in Collaborative Environments
### 2.4.1 Data Sharing
Ensuring privacy while allowing data sharing across multiple organizations is a significant challenge. Blockchain's transparency, where all transactions are visible to all participants, can conflict with the need for privacy, especially when sensitive data is involved.

### 2.4.2 Anonymization
Integrating privacy-preserving techniques, such as differential privacy or data anonymization, within a blockchain framework is complex and requires careful balancing to ensure that the data remains useful while still protecting individual privacy.

## 2.5 Integration Complexity
### 2.5.1 Interoperability
Integrating blockchain with existing big data tools, such as Hadoop and Spark, presents significant technical challenges. These systems were not originally designed to work together, and ensuring seamless interoperability requires significant modifications to both the blockchain and big data frameworks.

### 2.5.2 Legacy Systems
Many organizations rely on legacy systems for their big data operations. Integrating blockchain into these existing infrastructures without disrupting operations or requiring a complete overhaul of the system is another challenge that current research has not fully addressed.

## 2.6 Resource Intensity
### 2.6.1 Costs
The computational power required to maintain a blockchain, particularly in public networks, is substantial. When integrated with big data platforms, the resource demands increase further, as both systems require significant processing power and storage capacity. This can lead to higher operational costs and pose a barrier to widespread adoption.

Addressing these challenges is critical for realizing the full benefits of integrating blockchain technology with big data platforms. Successfully overcoming these issues will enable the creation of secure, scalable, and efficient data sharing systems that can be widely adopted across industries such as healthcare, finance, and supply chain management. Solving these problems will facilitate the development of robust, future-proof data ecosystems that combine the strengths of both blockchain and big data technologies, ultimately driving innovation and improving data management practices on a global scale.

## 3. Proposed Solution
The proposed solution involves the integration of two primary technologies: **Big Data Processing** using Apache Spark and **Blockchain Technology** for secure and immutable data logging. The integration aims to leverage the strengths of both technologies, addressing the challenges of scalability, data security, and real-time processing in high-throughput environments.

### 3.1 Overview
### 3.1.1 Big Data Processing with Apache Spark
Apache Spark is employed to handle large-scale datasets with high speed and efficiency. Spark's distributed processing framework allows for the parallel processing of data across a cluster of machines, making it ideal for handling big data operations that require rapid computation and scalability.

### 3.1.2 Blockchain Technology for Secure Logging
Blockchain is used to ensure the security, integrity, and immutability of data. By logging processed data onto a blockchain, the system guarantees that data records cannot be tampered with or altered, providing a reliable audit trail. The use of multiple blockchain nodes ensures redundancy and enhances the system's resilience.

## 3.2 Key Components
### 3.2.1 Data Loading and Preprocessing (Big Data Component):
**Objective:** Load and preprocess large datasets to prepare them for analysis and subsequent blockchain logging.

Process: Data is ingested into the Apache Spark framework, where it undergoes cleansing, transfor- mation, and aggregation. The data is then partitioned for parallel processing across the Spark cluster.

### 3.2.2 Parallel Processing and Blockchain Logging (Hybrid Component)
**Objective:** Process data in parallel using Apache Spark and log the processed data onto blockchain nodes.

**Process:** The preprocessed data is divided into batches, which are processed in parallel. Each batch's results are logged onto multiple blockchain nodes, ensuring that the data is securely stored and immutable.

### 3.2.3 Encryption and Key Management (Blockchain Component):
**Objective:** Secure the data before logging it onto the blockchain to protect sensitive information.

**Process:** Data is encrypted using advanced encryption techniques, and encryption keys are managed in a decentralized manner to ensure security and scalability.

### 3.2.4 Differential Privacy and Data Anonymization (Big Data Component)
**Objective:** Apply differential privacy techniques to protect sensitive data while maintaining its utility for analysis.

**Process:** Controlled noise is added to the data to anonymize it, preventing the re-identification of in- dividuals while allowing meaningful insights to be drawn from the data.

### 3.2.5 Performance Testing and Security Verification
**Objective:** Test and verify the performance and security of the integrated system.

**Process:** The system undergoes rigorous performance testing to ensure it can handle large volumes of data and high transaction throughput. Security tests verify the effectiveness of encryption, privacy measures, and the immutability of blockchain logs.

This hybrid architecture offers a robust solution for secure, scalable data processing and sharing across organizations. It is particularly well-suited for high-throughput environments where both the integrity and privacy of data are paramount. By combining the strengths of big data technologies and blockchain, the solution addresses the limitations of each, providing a comprehensive framework for modern data-driven applications.

## 4. Methodology
The proposed methodology integrates blockchain technology with big data processing frameworks to create a secure, scalable, and efficient system for managing large datasets. This hybrid architecture leverages the distributed computing power of Apache Spark for data processing and the immutability of Ethereum blockchain for secure data logging. The methodology is structured into several key components, each addressing specific aspects of the system's functionality, from data preprocessing to secure logging and performance evaluation.

## 4.1 Blockchain Setup and Deployment
**Objective:** Deploy and configure smart contracts on multiple blockchain nodes to establish a decentralized, secure logging mechanism.

**Mathematical Formulation:**
Let $B$ represent the set of blockchain nodes:

$$\mathcal{B} = \{B_1, B_2, \ldots, B_n\}$$

where each $B_i$ is a blockchain node configured with a smart contract $C$ defined by its ABI (Application Binary Interface) and bytecode. The deployment function $D(C, B_i)$ deploys the contract on node $B_i$:

$$D(C, B_i) \rightarrow \text{Contract Address}$$

This ensures that C is replicated across all nodes Bi, establishing a distributed logging system.

## 4.2 Data Loading and Preprocessing (Big Data Component)
**Objective:** Efficiently load, preprocess, and prepare large- scale datasets for distributed processing using Apache Spark, followed by secure logging on the blockchain.

Let $\mathbf{D} = \{d_1, d_2, \ldots, d_m\}$ represent the raw dataset where each $d_j$ is a data point.

• **Data Ingestion:** The raw dataset $\mathbf{D}$ is loaded into Spark's distributed Data Frame $\mathbf{F}$:

$\mathbf{F} = \text{Load Data}(\mathbf{D})$

• **Preprocessing:** Transform the dataset by applying cleansing operations $\mathcal{T}$:

$$\mathbf{F}' = \mathcal{T}(\mathbf{F})$$

where F′ is the cleaned and preprocessed Data Frame.

• **Partitioning:** Partition F′ into subsets $\mathbf{P}_k$ for parallel processing:

$\mathbf{P} = \{\mathbf{P}_1, \mathbf{P}_2, \ldots, \mathbf{P}_q\}$

where each partition $\mathbf{P}_k$ is processed independently across Spark's cluster nodes.

## 4.3 Parallel Data Processing and Blockchain Logging (Hybrid Component)
**Objective:** Process data in parallel using Apache Spark and securely log the processed results onto multiple blockchain nodes.

**Mathematical Formulation:**
̃ **Parallel Processing:** Each partition $\mathbf{P}_k$ is processed in parallel using a function $\mathcal{P}(\mathbf{P}_k)$:

$\mathbf{R}_k = \mathcal{P}(\mathbf{P}_k) \qquad (0, \frac{1}{\epsilon})$.

where $\mathbf{R}_k$ is the result of processing partition $\mathbf{P}_k$.

• **Blockchain Logging:** Log the result $\mathbf{R}_k$ onto the blockchain node $B_i$ using the contract $C$:

$\mathcal{L}(\mathbf{R}_k, B_i) \rightarrow \text{Immutable Log}$

This ensures that the processed data is securely stored in an immutable ledger.

## 4.4 Differential Privacy
**Objective:** Protect sensitive data by applying differential privacy

techniques, ensuring that the privacy of individual data points is preserved while maintaining the utility of the dataset for analysis.
Mathematical Formulation:
• **Epsilon-Differential Privacy:** The parameter $\epsilon$ controls the trade-off between privacy and accuracy. A smaller $\epsilon$ provides stronger privacy but introduces more noise into the data.
• **Laplace Mechanism:** Differential privacy is achieved by adding noise from the Laplace distribution to the dataset. Given a dataset **D** and a function $f(\mathbf{D})$, the Laplace mechanism adds noise $\eta$ to $f(\mathbf{D})$:
$f(\mathbf{D}) = f(\mathbf{D}) + \eta$

where $\eta \sim$ Laplace

**Implementation:**
• The dataset is first converted into a format suitable for numerical processing, specifically a NumPy ar-ray. This step facilitates efficient manipulation and application of mathematical operations necessary for differential privacy.
• Next, noise is generated and added to the dataset to ensure differential privacy. This noise is drawn from a Laplace distribution, where the scale parameter is inversely proportional to the privacy parameter $\epsilon$. The calculation follows:
noisy data = data + Laplace $(0, \frac{1}{\epsilon})$.

This addition of noise helps to obscure individual data points, thereby protecting privacy.
• Finally, the resulting dataset, now containing noisy data, is organized and prepared for subsequent processing steps. This preparation ensures that while privacy is protected, the dataset remains useful for analysis.

## 4.5 E. Key Management
**Objective:** Securely manage and distribute cryptographic keys used for encryption and decryption in a decentralized environment, ensuring data confidentiality.
**Mathematical Formulation:**
**Key Generation:** A unique cryptographic key K is generated by hashing the concatenation of an identity ID and a set of attributes Attr:
$K = \text{SHA256}(\text{ID} \,\|\, \text{Attr})$
Here, ID is the identity (e.g., a user ID) and Attr represents the associated attributes (e.g., roles or permissions).
**Encryption:** The plaintext message M is encrypted using AES (Advanced Encryption Standard) in CBC (Cipher Block Chaining) mode with the key $K$:
$C = \text{AES-CBC}(K, M)$
The ciphertext C is a combination of the initialization vector (IV) and the encrypted message.
**Decryption:** The ciphertext $C$ is decrypted using the same key $K$, and the original plaintext $M$ is retrieved:
$M = \text{AES-CBC}^{-1}(K, C)$
This ensures that only entities with the correct key can access the decrypted data.

## 4.6 Off-Chain Data Management
**Objective:** Efficiently manage large datasets by storing them off-chain while ensuring their integrity and availability, reducing the load on the blockchain.
**Mathematical Formulation:**
**Data Storage**: Data $D$ is stored off-chain using a decentralized storage system, such as IPFS (Inter Planetary File System). The process involves first encoding the data and then generating a unique identifier for it, known as the IPFS hash IPFS hash. Additionally, a cryptographic hash $H(D)$ is computed using the SHA-256 algorithm to ensure data integrity:
IPFS hash = unique identifier from IPFS($D$)
$H(D) = \text{SHA256}(D)$
The process returns both the IPFS hash, which serves as the address for retrieving the data, and the cryptographic hash, which can be used later for verifying the integrity of the data.
**Data Verification:** To verify the integrity of the data, the stored hash $H(D)$ is compared with a newly computed hash $H'(D)$ of the retrieved data:

$$\text{Verify}(D, H(D)) = \begin{cases} \text{True} & \text{if } H(D) = H'(D) \\ \text{False} & \text{otherwise} \end{cases}$$

This process ensures that the data has not been tampered with and is identical to the original stored data.

## 4.7 Performance Testing and Security Verification
**Objective:** Evaluate the system's performance under high data loads and verify the security of the encryption and logging mechanisms.
**Mathematical Formulation:**
**Scalability Testing:** Measure the system's performance as a function of the scale factor s:
Performance($s$) = f($s$)
where s represents the dataset size or transaction volume.
**Security Testing:** The effectiveness of encryption and differential privacy is tested by evaluating the probability of data breach Pr(Breach), which should be minimized:
Pr(Breach) $\approx 0$

## 4.8 Report Generation and Monitoring
**Objective:** Compile and report the system's performance and security metrics.
**Mathematical Formulation:**
**Performance Report:** Summarize the results of scalability and processing times $T(s)$:
Report = {Performance($s$), $T(s)$}
**Security Report:** Document the results of encryption and privacy tests, ensuring compliance with security standards.
This comprehensive methodology integrates advanced big data processing capabilities with the security and immutability of blockchain technology. By leveraging differential privacy, cryptographic key management, and off-chain data storage, the system provides robust, scalable, and secure data management.

## 5. Experiments and Evaluation
### 5.1 Experiment Setup
The experiments were conducted to evaluate the performance, scalability, and security of the proposed hybrid system that integrates blockchain with big data processing. The setup involved two distinct experiments, followed by performance and security testing.

### 5.1.1 Technologies Used
• **Apache Spark:** Used for distributed big data processing, specifically for data loading, cleansing,
transformation, and aggregation.
• **Ethereum Blockchain:** Deployed with Vyper smart contracts across multiple nodes to log data securely and immutably.
• **IPFS (Inter Planetary File System):** Utilized for off-chain data storage, reducing the burden on the blockchain by storing large datasets off-chain.
• **IABHE (Identity and Attribute-Based Honey Encryption):** Applied for encrypting sensitive data before it is logged onto the blockchain.
• **Differential Privacy:** Ensured data privacy by adding Laplace noise to the dataset.

### 5.1.2 System Configuration
• **Single Blockchain Server Setup:** In the first experiment, a single blockchain node was deployed, and linear insertion of data was performed to measure end-to-end performance time.
• **Hybrid Setup with 5 Blockchain Servers:** In the second experiment, five blockchain nodes were deployed to parallelize the blockchain processing. Apache Spark was used for big data processing to compare the time consumption and performance against the single-node setup.

### 5.2 Experiment Execution
### 5.2.1 Single Blockchain Server Experiment
**Objective:** Measure the performance of a single blockchain node handling linear data insertion.
**Process:**
*Data was loaded and preprocessed using Apache Spark.
∗ The preprocessed data was sequentially inserted into the single blockchain node, and the total time
taken for this process was recorded.

### 5.2.2 Hybrid Setup with 5 Blockchain Servers and Apache Spark
**Objective:** Evaluate the performance and scalability of a hybrid setup with parallel blockchain processing and big data handling using Spark.
**Process:**
∗ The same data was loaded and preprocessed using Apache Spark as in the first experiment.
∗ The preprocessed data was chunked based on timestamps, with each chunk containing multiple records indexed by their transaction time. These chunks were distributed across the 5 blockchain nodes for parallel processing.
*Apache Spark was used to handle the data processing, and the

blockchain nodes logged the data
concurrently.
∗ The total time taken for the entire process was compared against the single-node setup to assess
improvements in performance and scalability.

### 5.3 Performance Testing
Performance testing was conducted using the following methods:

### 5.3.1 Data Processing Time Measurement
Method: The time taken to process the data within Apache Spark was measured by increasing the
dataset size incrementally (simulating scalability) and recording the processing time. For load testing,
the data size was doubled in each iteration, with a total of 10 iterations, and the big data processing
time was gathered.
**Blockchain Logging Scalability:**
**Method:** The time taken to log data onto the blockchain was measured by simulating an increase
in data size (scale factor). The experiment was conducted for both the single blockchain node and the
5-node setup.
**Chunked Data Processing:**
**Method:** The data was chunked based on timestamps, with each transaction time indexing multiple
records. This method allowed for efficient batch processing and logging, reducing the overall transaction
time.

### 5.3.2 Security Testing
Security testing focused on evaluating the encryption and differential privacy mechanisms:
**Encryption and Decryption Testing:**
**Method:** Data was encrypted using IABHE and then decrypted to verify the integrity of the process.
The test ensured that the original plaintext data was accurately recovered after decryption.
**Differential Privacy Testing:**
**Method:** Laplace noise was added to the dataset to ensure differential privacy. The test verified that the
noisy data differed from the original data and that the noise was sufficient to protect individual data points while maintaining overall data utility.

## 6. Results
### 6.1 Scalability Testing
The performance of the big data processing component was evaluated by progressively increasing the size of the dataset and measuring the corresponding scalability time. The dataset size was doubled in each iteration, starting from 390,424 records, and the scalability time was recorded. Below is the graph depicting the relationship between the number of records processed and the time taken for processing.
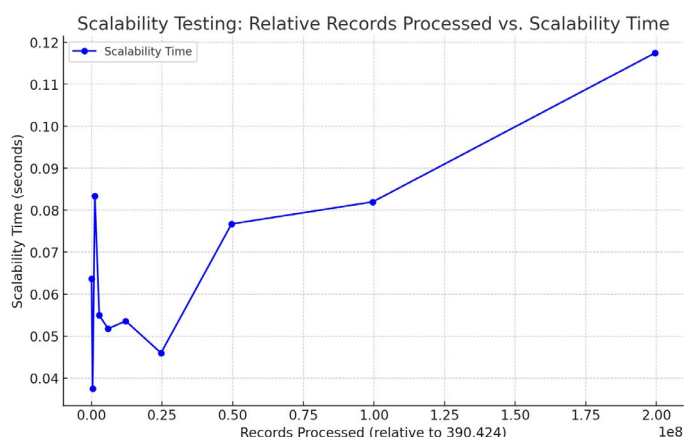
**Figure 1:** Scalability Testing: Records Processed vs. Scalability Time

## 6.2 Metrics Comparison

The following table summarizes the key metrics observed during the experiments:

| Metric | Single Node Setup | Hybrid Setup (5 Nodes) |
|---|---|---|
| Records Processed | 390,424 | 390,424 |
| Overall Execution Time (seconds) | 16026.21 | 3009.81 |
| Encryption/Decryption Test | Success | Success |
| Differential Privacy Test | Success | Success |

**Table I: Comparison of Key Performance Metrics Between Single Node and Hybrid Setup**

## 7. Discussion

The comparative analysis between the single-node and hybrid setups reveals significant performance differences, particularly in terms of overall execution time, scalability, and efficiency.

## 7.1 Overall Execution Time

The overall execution time for the single-node setup was 16,026.21 seconds, compared to 3,009.81 seconds for the hybrid setup. This represents a **reduction of approximately 81.22%** in execution time when using the hybrid model. The significant decrease in time underscores the efficiency gains achieved through parallel processing across 5 blockchain nodes. By distributing the workload, the hybrid setup was able to process the data pipeline much faster, overcoming the bottlenecks observed in the single-node configuration.

## 7.2 Scalability and Processing Time

The scalability tests further illustrate the system's ability to handle increasing data volumes effectively. The tests involved doubling the dataset size in each iteration, starting from 390,424 records and scaling up to nearly 200 million records. The hybrid setup demonstrated impressive scalability:

**Initial Processing:** For the first set of 390,424 records, the system took approximately 0.064 seconds.

**Maximum Processing:** When the number of records increased to 199,897,088, the processing time only in- creased to 0.117 seconds.

This indicates that despite a **51-fold increase** in the number of records, the processing time only increased by **83%** (from 0.064 seconds to 0.117 seconds). In contrast, the single-node setup, with a total execution time of 16,026.21 seconds, would have struggled to maintain such performance levels with in- creasing data volumes, likely leading to exponential increases in processing time. The hybrid model's ability to distribute the processing workload across multiple nodes allowed it to maintain a nearly linear relationship between the number of records processed and the time taken, ensuring that the system did not become overwhelmed as data volumes increased.
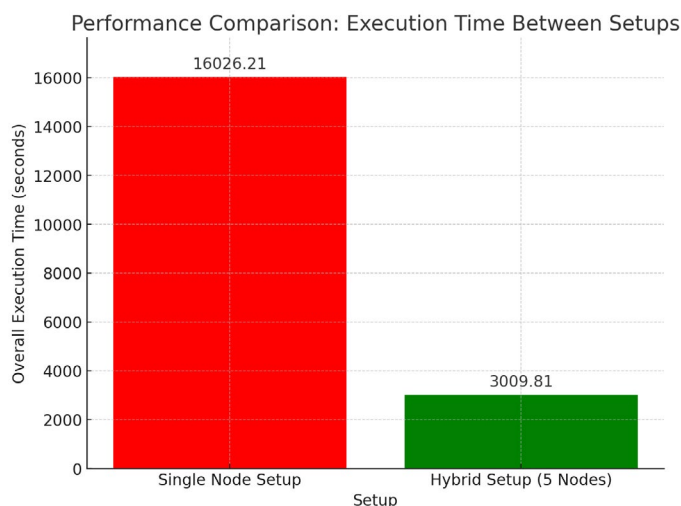


**Figure 2:** Performance Comparison: Execution Time Between Single Node and Hybrid Setup (5 Nodes)

## 7.3 Speedup Analysis

To quantify the efficiency gains, the speedup achieved by the hybrid setup can be calculated as:

$$\text{Speedup} = \left( \frac{\text{Execution Time (Single Node)}}{\text{Execution Time (Hybrid)}} \right) \approx 5.33$$

This means that the hybrid setup processes the same work- load approximately 5.33 times faster than the single-node setup. Furthermore, the **Speedup Percentage** can be calculated as:

$$\text{Speedup Percentage} = \left( \frac{16026.21 - 3009.81}{16026.21} \right) \times 100 \approx 81.22\%$$

this indicates an **81.22% reduction** in execution time, clearly demonstrating the significant efficiency improvements offered by the hybrid setup.

## 7.4 Security and Data Integrity

Both setups successfully passed the encryption/decryption and differential privacy tests. This consistency across different configurations demonstrates that the system's security mechanisms were robust and reliable, regardless of the setup used. The encryption ensured that sensitive data was protected during

processing and storage, while differential privacy added a layer of protection against the re-identification of individual data points. These results confirm that the hybrid setup does not compromise on security while offering improved performance.

## 7.5 Overall System Performance

The hybrid setup not only reduced overall execution time by 81.22% compared to the single-node setup but also demonstrated a strong ability to scale efficiently with increasing data volumes. This performance improvement, combined with robust security features, makes the hybrid model an optimal solution for large-scale, high-throughput environments. Organizations dealing with extensive and sensitive data can greatly benefit from the hybrid approach, as it ensures both efficiency and security in a decentralized processing environment.

In conclusion, the hybrid blockchain-big data integration model offers significant advantages over a traditional single- node setup, particularly in terms of reduced execution time, better scalability, and maintained security. These benefits position the hybrid approach as a superior solution for managing large datasets in distributed and decentralized systems.

## 7.6 Cost Comparison

The cost comparison between the single-node setup and the hybrid setup with 5 nodes is based on the assumption that the cost of running a blockchain node is $0.0001 per second per node. This fixed cost is applied to both setups to determine the total cost of execution.

*Single-Node Setup:* For the single-node setup, the total execution time was 16,026.21 seconds. The total cost is calculated as:

$$\text{Cost}_{\text{Single Node}} = \text{Execution Time} \times \text{Cost per Second per Node}$$

$$\text{Cost}_{\text{Single Node}} = 16{,}026.21 \times 0.0001 = \$1.6026$$

*Hybrid Setup (5 Nodes):* For the hybrid setup, with 5 nodes operating in parallel, the total execution time was significantly reduced to 3,009.81 seconds. However, the cost per second is multiplied by the number of nodes (5 nodes):

$$\text{Cost}_{\text{Hybrid}} = \text{Execution Time} \times \text{Cost per Second per Node} \times \text{Number of Nodes}$$

$$\text{Cost}_{\text{Hybrid}} = 3{,}009.81 \times 0.0001 \times 5 = \$1.5049$$

*Cost Efficiency:* Despite the hybrid setup involving multiple nodes, it remains slightly less expensive than the single- node setup due to the significant reduction in execution time. Specifically, the hybrid setup costs $1.5049, while the single- node setup costs $1.6026. This results in a small cost saving of approximately 6.1%, demonstrating that the hybrid setup is not only more time-efficient but also more cost-effective under the assumed cost conditions.

The cost analysis indicates that the hybrid setup, with its parallel processing capability, provides both time and cost savings compared to the single-node setup. Even though the hybrid setup uses more nodes, the reduction in execution time outweighs the additional cost of running multiple nodes, making it a superior choice for both performance and cost efficiency.

## 8. Future Scope

The integration of blockchain technology with big data platforms presents numerous opportunities for future research and development. As this field continues to evolve, several avenues can be explored to enhance the effectiveness, scalability, and applicability of hybrid blockchain-big data systems. Below are key areas for future work:

### 8.1 Enhanced Scalability Solutions

While the current hybrid setup demonstrates improved scalability, future research could focus on developing more advanced techniques to further enhance the scalability of such systems. This could involve exploring sharing methods in blockchain technology, optimizing data partitioning strategies in big data platforms, or integrating new consensus algorithms tailored for high-throughput environments.

### 8.2 Energy Efficiency and Sustainability

Blockchain operations, especially in a distributed multi- node environment, can be energy-intensive. Future work could investigate ways to reduce the energy footprint of hybrid systems. This might include developing energy-efficient consensus mechanisms, utilizing green energy sources for blockchain nodes, or optimizing the data processing pipeline to minimize unnecessary computational overhead.

### 8.3 Integration with Emerging Technologies

The hybrid architecture could be further enriched by integrating with emerging technologies such as artificial intelligence (AI), machine learning (ML), and the Internet of Things (IoT). For instance, AI/ML algorithms could be used to predict and manage workloads more efficiently, while IoT devices could feed real-time data into the blockchain, enhancing the system's responsiveness and adaptability.

### 8.4 Security Enhancements

Although the current system includes encryption and differential privacy to ensure data security and privacy, future work could explore more robust security frameworks. This could involve implementing post-quantum cryptography to safeguard against future quantum computing threats, or developing more sophisticated privacy-preserving techniques that enable secure data sharing without compromising individual privacy.

### 8.5 Cross-Chain Interoperability

As multiple blockchain networks become more prevalent, enabling cross-chain interoperability will be crucial. Future re- search could focus on creating frameworks that allow seamless interaction between different blockchain networks, enabling data and assets to move freely and securely across chains without compromising the integrity or security of the data.

## 8.6 Real-Time Data Processing
Enhancing the real-time data processing capabilities of the hybrid system could be another area of exploration. This could involve integrating stream processing frameworks with the blockchain to handle high-velocity data, ensuring that the system can process and log data as it is generated with minimal latency.

## 8.7 Decentralized Data Marketplaces
The development of decentralized data marketplaces, where data providers and consumers can interact securely and transparently using blockchain technology, is a promising area for future work. Such marketplaces could leverage the hybrid architecture to ensure data integrity, privacy, and secure trans- actions, fostering trust in data exchange.

## 8.8 Regulatory Compliance and Governance
As regulatory environments evolve, particularly concerning data privacy (e.g., GDPR), the hybrid system could be adapted to ensure compliance with various international regulations. Future research could focus on developing governance frame- works that integrate regulatory compliance into the blockchain and big data architecture, ensuring that the system adheres to legal standards while maintaining operational efficiency.

## 8.9 Usability and Adoption
Finally, increasing the usability and adoption of hybrid blockchain-big data systems in industry and government sec- tors is a critical future direction. This could involve developing user-friendly interfaces, offering modular solutions that can be easily integrated into existing infrastructures, and conducting case studies or pilot projects to demonstrate the practical benefits of the system in real-world scenarios.

The hybrid blockchain-big data integration model holds significant potential for addressing some of the most pressing challenges in secure, scalable data management. By exploring these future directions, researchers and practitioners can further enhance the capabilities and applications of this technology, making it a cornerstone of next-generation digital infrastructures.

## References
1. Siyal, R., Long, J., Asim, M., Ahmad, N., Fathi, H., & Alshinwan, M. (2024). Blockchain-Enabled Secure Data Sharing with Honey Encryption and DSNN-Based Key Generation. *Mathematics, 12*(13), 1956.
2. Tekchandani, P., Bisht, A., Das, A. K., Kumar, N., Karuppiah, M., Vijayakumar, P., & Park, Y. (2024). Blockchain-Enabled Secure Collaborative Model Learning using Differential Privacy for IoT-Based Big Data Analytics. *IEEE Transactions on Big Data*.
3. Ren, W., Zhang, W., Liu, J., Cai, H., & Liu, H. (2023, October). Blockchain-Based Data Security Sharing System. In *Proceedings of the 2023 7th International Conference on Electronic Information Technology and Computer Engineering* (pp. 1006-1009).
4. Choubey, A., Choubey, S., Jaiswal, D., & Jaiswal, M. (2024, March). Integrating Blockchain in Cloud Computing for Enhanced Data Management and Security. In *2024 11th International Conference on Reliability, Infocom Technologies and Optimization (Trends and Future Directions)(ICRITO)* (pp. 1-7). IEEE.
5. Zou, Y., Peng, T., Wang, G., Luo, E., & Xiong, J. (2023). Blockchain-assisted multi-keyword fuzzy search encryption for secure data sharing. *Journal of Systems Architecture, 144*, 102984.