# Integrating K-Means Clustering with PoA Blockchain and IPFS for Clustered Data Synchronization: The OriBloX CDSF Approach

Edy Saputro
Master of Information System, Postgraduate School, Diponegoro University, Indonesia

Mustafid
Department of Statistics, Diponegoro University, Semarang, Indonesia

Jatmiko Endro Suseno
Department of Physics, Diponegoro University, Semarang, Indonesia

## ABSTRACT

Efficient data synchronization in distributed systems presents significant challenges, as centralized solutions often face limitations in scalability, bandwidth efficiency, and resilience to single points of failure. Existing blockchain and decentralized storage technologies struggle to manage frequent data updates effectively. To address these issues, OriBloX CDSF integrates K-Means clustering (optimized with the Elbow method), TF-IDF analysis, Hyperledger Besu (using QBFT PoA consensus), and IPFS to deliver a secure, scalable, and decentralized synchronization framework. Its selective synchronization mechanism optimizes bandwidth usage by retrieving only updated cluster files, reducing unnecessary data transfers by up to 70%. Using the Amazon product catalog dataset, the framework demonstrated robust clustering performance, with the Elbow method consistently identifying optimal clusters and silhouette scores reaching up to 0.114, reflecting well-separated and meaningful groupings. OriBloX's design ensures efficient metadata synchronization, scalability, and fault tolerance, making it a reliable solution for distributed ecosystems.

## General Terms

Security, Algorithms, Data Synchronization, Distributed Systems, Blockchain Technology, Cloud Computing, Fault Tolerance, Decentralization.

## Keywords

K-Means Clustering, TF-IDF, Elbow Method, Blockchain, Proof of Authority, QBFT, Hyperledger Besu, IPFS, Data Synchronization, Distributed Systems, Decentralization, Fault Tolerance, Smart Contracts.

## 1. INTRODUCTION

### 1.1 Background

Modern data-driven ecosystems often operate across geographically distributed nodes, requiring timely and accurate access to shared information. Centralized architecture, while simpler to implement, faces significant challenges such as limited scalability, single points of failure, and inefficient bandwidth usage. These limitations become even more pronounced as datasets grow in complexity and size, particularly when handling dynamic and text-rich information such as product descriptions or customer reviews.

One promising solution to overcome these challenges is the use of blockchain technology. Blockchain offers a decentralized approach by providing an immutable ledger for metadata storage, supported by robust consensus mechanisms [10].

Unlike centralized systems, blockchain ensures data integrity and prevents single points of failure. However, while blockchain excels at securely storing metadata, it is impractical for large datasets due to high on-chain storage costs and limited performance scalability [10].

To address the issue of storing large datasets, decentralized file systems like IPFS (InterPlanetary File System) have been developed. IPFS complements blockchain by offering off-chain, content-addressable storage that reduces on-chain overhead [5]. However, while IPFS is efficient for storage, it lacks robust mechanisms to synchronize frequently updated or partially changed datasets across distributed nodes, which is a critical requirement in modern ecosystems.

Building on these advancements, OriBloX CDSF (Clustered Data Synchronization Framework) combines the strengths of blockchain, IPFS, and clustering techniques to create an integrated solution. By leveraging K-Means clustering (optimized via the Elbow method) and TF-IDF text analysis, OriBloX enables efficient organization of text-based datasets [3], [8]. This data is securely synchronized across nodes using Hyperledger Besu's permissioned blockchain with QBFT (Quorum Byzantine Fault Tolerance) consensus and stored off-chain using IPFS [4], [10]. Furthermore, OriBloX introduces a selective synchronization mechanism, allowing nodes to fetch only updated or relevant cluster files, thereby reducing bandwidth usage and improving efficiency. This hybrid approach offers a scalable, secure, and fault-tolerant solution for distributed data synchronization.

### 1.2 Problem statement

Modern distributed environments operate under heterogeneous conditions, with nodes varying in network capacities, storage resources, and processing power. Despite these variations, timely and consistent access to shared information remains critical. Traditional centralized solutions often lead to single points of failure and lack the scalability and fault tolerance required for data-intensive applications [9]. This issue is further exacerbated by the increasing volume and velocity of textual data, such as product descriptions and customer reviews, necessitating efficient clustering techniques for semantic grouping [8].

Blockchain technology offers a decentralized approach to secure metadata storage, yet it is unsuitable for storing large datasets due to excessive costs and performance constraints [10]. Similarly, while IPFS (InterPlanetary File System) enables scalable off-chain storage and decentralized data access, it lacks robust synchronization mechanisms for frequently updated or partially modified datasets [5].

The challenge becomes more complex in scenarios requiring selective updates, where nodes need only specific subsets of data, such as individual product clusters. Current systems force nodes to re-download entire datasets for minor changes, resulting in inefficient bandwidth usage and operational overhead. This limitation hinders scalability and responsiveness in real-time or near-real-time applications.

## 1.3 Objective

To address these challenges, this research proposes the OriBloX CDSF framework, which integrates K-Means clustering, TF-IDF analysis, Hyperledger Besu (using QBFT PoA consensus), and IPFS to provide a scalable and efficient solution. The research objectives are to:

1. Develop an efficient clustering mechanism for text-based datasets using K-Means and TF-IDF.

2. Securely manage clustering metadata through a permissioned blockchain using QBFT consensus.

3. Reduce on-chain storage costs by leveraging IPFS for scalable off-chain data distribution.

4. Implement a selective synchronization mechanism to minimize bandwidth usage by fetching only updated or relevant cluster files.

5. Ensure high fault tolerance and scalability, enabling robust performance even under intermittent node failures.

## 2. RELATED WORK

Efficient data management in distributed systems has led to significant advancements in blockchain and decentralized storage technologies, clustering algorithms, and hybrid frameworks. OriBloX CDSF builds on these innovations to address critical gaps in scalability, fault tolerance, and synchronization efficiency.

## 2.1 Blockchain and decentralized storage for data management

Blockchain technology, initially designed for decentralized financial systems [1], has evolved into enterprise-grade platforms like Hyperledger Besu, supporting permissioned environments with robust governance [4]. Proof of Authority (PoA) consensus mechanisms significantly reduce computational overhead compared to Proof of Work [13], [18], making them ideal for metadata synchronization in distributed systems.

To address the limitations of on-chain storage, the InterPlanetary File System (IPFS) provides decentralized, content-addressable storage, reducing blockchain storage overhead [5]. By enabling selective retrieval of data via Content Identifiers (CIDs), IPFS enhances scalability and bandwidth efficiency [11]. Additionally, its versioning and immutability features ensure a transparent update history, which is particularly useful for dynamic datasets. OriBloX CDSF combines PoA consensus with IPFS for secure, scalable metadata synchronization and efficient off-chain storage, addressing the limitations of existing solutions.

## 2.2 Clustering and synchronization techniques in distributed Systems

K-Means clustering, widely used for its simplicity and efficiency, enables the discovery of meaningful patterns in large datasets [3]. TF-IDF enhances clustering precision by emphasizing term importance and reducing noise from irrelevant words [8]. While distributed implementations of K-Means [12], [14] improve computational efficiency, they often rely on centralized coordination and lack mechanisms for secure synchronization of cluster outputs.

Synchronization remains a longstanding challenge in distributed systems [9]. Traditional replication protocols ensure consistency but fall short in optimizing bandwidth usage and handling dynamic datasets. Selective updates, which fetch only the changed portions of data, significantly reduce bandwidth requirements [16]. OriBloX CDSF addresses these challenges by integrating clustering with blockchain and IPFS, implementing selective synchronization mechanisms to fetch only updated cluster files. This approach ensures efficient, secure, and synchronized clustering in distributed ecosystems.

## 2.3 Gaps addressed by OriBloX CDSF

Despite advancements in blockchain, IPFS, and clustering, several critical gaps remain:

1. **Text-Focused Clustering**: Existing frameworks often prioritize numeric data types or supervised learning algorithms, leaving text-based clustering applications underexplored [8], [19]. This gap is particularly significant given the prevalence of text-rich datasets in industries such as e-commerce and social media analytics.

2. **End-to-End Framework**: Few solutions provide an integrated pipeline combining clustering, blockchain metadata storage, and off-chain data synchronization. For example, hybrid blockchain-IPFS approaches [17], [19] highlight the potential but fail to address selective synchronization and clustering for dynamic datasets.

3. **Robust Fault Tolerance**: Research on fault tolerance in permissioned blockchain environments is limited. Most studies focus on public blockchain fault tolerance [1], [13], neglecting the unique challenges posed by frequent updates in distributed clusters.

OriBloX CDSF addresses these gaps by integrating K-Means clustering, TF-IDF analysis, PoA consensus, and IPFS with selective synchronization. This unified framework ensures scalable, secure, and efficient data management in distributed ecosystems, overcoming the limitations of existing solutions [5], [15].

## 3. METHODOLOGY

## 3.1 OriBloX CDSF design and workflow

OriBloX CDSF integrates K-Means clustering, TF-IDF analysis, Hyperledger Besu (using QBFT PoA consensus), and IPFS to enable secure, scalable, and decentralized data synchronization. The workflow involves several stages:

1. **Data Preprocessing**: Cleaning and normalizing raw datasets, converting text into TF-IDF feature vectors.

2. **Clustering**: Using the Elbow Method to optimize cluster formation and generating unique Cluster IDs with SHA-256.

3. **Metadata Management**: Updating cluster metadata on-chain using smart contracts for secure synchronization.

4. **Selective Synchronization**: Fetching only updated data via IPFS, reducing bandwidth usage and ensuring consistent synchronization.
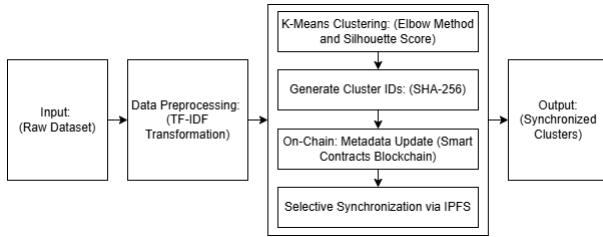
**Fig 1: OriBloX CDSF workflow diagram**

## 3.2 Data preprocessing and feature extraction

The Amazon product catalog dataset from Kaggle, containing attributes such as titles, descriptions, and categories, serves as the input. Data preprocessing involves cleaning missing entries, normalizing text fields, and scaling numerical attributes for uniformity. Text preprocessing includes tokenization, stopword removal, and stemming to standardize terms. These processes are executed by Electron Nodes, ensuring consistent data preparation across nodes. TF-IDF transformation converts textual data into feature vectors, with optional dimensionality reduction using PCA for computational efficiency.

The Amazon product catalog dataset from Kaggle was selected for this study due to its large volume of textual data and diverse product categories, making it an ideal benchmark for evaluating text-based clustering approaches. This dataset includes rich semantic attributes (e.g., product titles, descriptions, and categories), which align well with TF-IDF-based feature extraction and K-Means clustering techniques. Additionally, this dataset represents a real-world e-commerce environment, where efficient data synchronization and distributed clustering are critical for applications like product recommendations, inventory management, and personalized search results.

Given the complexity and scalability of the dataset, OriBloX CDSF's performance on this benchmark provides robust evidence of its applicability in real-world large-scale distributed environments. While alternative datasets could further validate the model's generalizability, the current evaluation focuses on demonstrating core clustering efficiency and synchronization performance.

## 3.3 Clustering process

K-Means clustering optimizes cluster formation using the Elbow Method to evaluate WCSS and determine the optimal k value. Centroids are initialized using k-means++ and iteratively refined until convergence. Cluster IDs are generated by hashing centroid vectors (using SHA-256) and are used to track clusters across IPFS and QBFT. Each cluster is saved as a separate JSON file (e.g., cluster_shirt.json) for efficient management and retrieval.
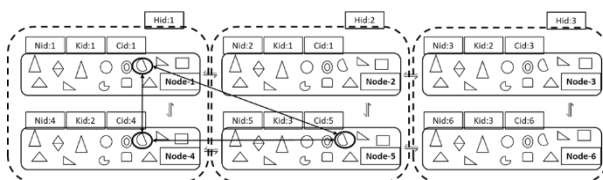


**Fig 2: OriBloX CDSF cluster id generation**

- Compute centroid vectors for each cluster.
- Hash these vectors use SHA-256 to generate unique Cluster IDs, for tracking clusters across IPFS and QBFT.

- After generating clusters, save each cluster to a separate file:
    - cluster_shirt.json
    - cluster_jacket.json
    - cluster_pants.json
- Each file represents a unique Cluster ID.

## 3.4 Off-Chain storage via IPFS

OriBloX CDSF leverages IPFS (InterPlanetary File System) to store cluster data off-chain, ensuring decentralized and tamper-resistant storage. Each cluster file (e.g., cluster_shirt.json) contains product IDs, centroids, and summary statistics. Instead of storing large data directly on the blockchain, IPFS generates immutable Content Identifiers (CIDs), which serve as unique references to each stored cluster.

### 3.4.1 Storage and Retrieval Process
1. Cluster File Creation
   - After clustering, each dataset is stored in a JSON format.
   - These files contain clustered product data, metadata, and centroids.
2. IPFS CID Generation
   - Each file is added to IPFS using the ipfs add command.
   - IPFS assigns a unique CID, representing the file's content hash using:

        ipfs add cluster_shirt.json

        Output: QmXkY...8H3

3. Decentralized Retrieval
   - Nodes retrieve cluster files using:

        ipfs get <CID>

   - Since IPFS follows content-based addressing, the same CID always retrieves the exact file.

### 3.4.2 Handling Data Updates
1. New Cluster Versions
   - When cluster data is updated, a new CID is generated since IPFS assigns hashes based on content.

        ipfs add cluster_shirt_updated.json

        Output: QmNewUpdatedCID

   - The updated CID replaces the previous one in the blockchain.
2. Version Control
   - Old versions remain accessible unless they are removed from the IPFS network.
   - OriBloX nodes can fetch previous cluster states if needed.

## 3.5 Metadata management and selective synchronization

OriBloX smart contracts manage metadata for clusters by mapping Cluster IDs to CIDs and associated metadata (e.g., timestamps, origin nodes). Metadata updates are validated using QBFT consensus, ensuring tamper-proof records. Selective synchronization involves:

1. Comparing local and on-chain CIDs to identify updated files.

2. Fetching updated data via IPFS using ipfs get <CID>.

3. Replacing or merging files locally, ensuring consistent cluster updates across nodes.

To illustrate metadata updates, below is a pseudocode representation of the OriBloX smart contract operations:

```
// OriBloX Smart Contract for Metadata Management

contract OriBloXCluster {
    mapping(string => string) private clusterCID;
    mapping(string => string) private metadata;

    function updateCluster(string memory clusterId, string
memory cid, string memory meta) public {
        clusterCID[clusterId] = cid;
        metadata[clusterId] = meta;
    }

    function getClusterCID(string  memory  clusterId)
public view returns (string memory) {
        return clusterCID[clusterId];
    }

    function getClusterMetadata(string memory clusterId)
public view returns (string memory) {
        return metadata[clusterId];
    }
}
```

This contract enables secure metadata registration, ensuring only valid updates are stored on-chain. The selective synchronization process ensures that nodes only fetch cluster files with updated CIDs, reducing unnecessary data transfers and improving bandwidth efficiency.

## 3.6 Deployment design

OriBloX CDSF operates within a distributed ecosystem consisting of:

1. **OriBloX Node (Electron):** Deployed on Dockerized virtual machines, handling local IPFS storage, QBFT synchronization, data ingestion, and clustering.

2. **Neutron Clusters:** Single-provider clusters deployed in a controlled laboratory environment, optimized for low-latency, enterprise use.

3. **Proton Hub:** Designed for future cross-provider collaboration but excluded from this research scope.
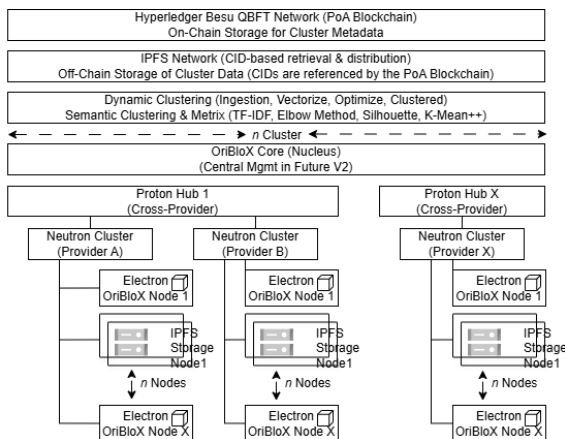


**Fig 3: OriBloX CDSF architecture**

OriBloX CDSF combines advanced clustering techniques, blockchain metadata management, and decentralized storage to provide a scalable and efficient framework for distributed data synchronization.

# 4. RESULTS
## 4.1 Clustering performance
### 4.1.1 Dataset & Elbow method
The OriBloX CDSF clustering was evaluated using the Amazon product catalog dataset from Kaggle. The original dataset was split into smaller, manageable parts (e.g., amazon_part_1, amazon_part_2, amazon_part_3, etc.) for efficient processing. Each part was treated as an independent dataset for analysis. For each dataset (e.g., Dataset1, Dataset2, Dataset3, etc.), the Elbow Method was applied to identify the optimal number of clusters ($k$). Key observations:

- **Dataset 1**: The WCSS plot exhibited a clear bend at k=8, suggesting an optimal cluster count of k=8.
- **Dataset 2:** The WCSS plot showed the optimal k=10, reflecting a strong bend point.
- **Dataset 3:** The WCSS plot indicated the optimal k=9, demonstrating another distinct bend point.

**Table 1. WCCS analysis**

| $k$ | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
| **Dataset 1 (WCSS)** | 5800 | 5700 | 5601 | 5500 | 5451 | 5401 | 5351 | 5301 | 5250 |
| **Dataset 2 (WCSS)** | 5800 | 5705 | 5601 | 5501 | 5450 | 5401 | 5350 | 5300 | 5251 |
| **Dataset 3 (WCSS)** | 5800 | 5700 | 5605 | 5500 | 5451 | 5400 | 5350 | 5301 | 5250 |

The Amazon product catalog dataset was selected for this study due to its large-scale textual data, diverse categories, and real-world relevance in distributed environments. The dataset contains rich semantic features (e.g., product descriptions, titles, and categories), making it ideal for TF-IDF-based clustering. This mirrors real-world applications such as recommendation systems, product classification, and federated learning.

Additionally, the dataset was split into smaller subsets (Dataset1, Dataset2, Dataset3) to simulate distributed clustering across multiple OriBloX nodes, ensuring that the evaluation reflects real-world decentralized clustering scenarios.

### 4.1.2 Silhouette score
The Silhouette Score was calculated for each dataset to evaluate clustering effectiveness. The score measures how well-separated the clusters are, with values closer to 1 indicating strong separation and values close to 0 suggesting overlapping clusters.

- **Dataset 1**: At the optimal k=8, the average silhouette score was **0.061**, indicating moderately well-separated clusters with limited overlap in TF-IDF space
- **Dataset 2**: At the optimal $k$=10 (as identified by the Elbow Method and Silhouette analysis), the average silhouette score was **0.109**, indicating moderately well-separated clusters with limited overlap in TF-IDF space.
- **Dataset 3**: At the optimal $k$=9 (as determined through optimization analysis), the silhouette score was **0.114**, reflecting better-separated clusters compared to Dataset 2.

These results confirm that higher k-values (Dataset 2 & 3) resulted in better-defined clusters, suggesting that TF-IDF+K-Means effectively captures product category distinctions.

### 4.1.3 Cluster composition

Clusters showed semantic coherence. Cluster IDs such as c63986122a19ed35ca51de317ff3a3c6c96657b4fc4c3fe9cc11f bb0e84a13ef grouped in corresponding categories (e.g., "Set" products of different sizes such as S, M, and XXL).

Observations from multiple datasets revealed that items like "men's sneakers" consistently grouped under distinct clusters, demonstrating the framework ability to segregate product categories effectively.
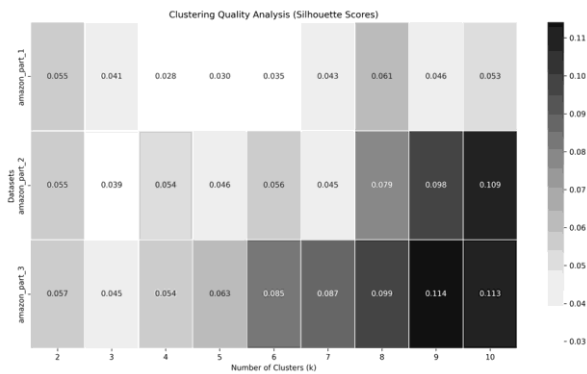


**Fig 4: Clustering quality heatmap**

### 4.1.4 Centroid-based ID

Each cluster was represented by its centroid vector, hashed using SHA-256 to generate unique Cluster IDs, ensuring traceability across nodes in the OriBloX network. For instance: Dataset 3's metadata highlights Cluster ID 934af5075b0f97aa1bc42c826f6aae872b361a1f090ff501f530b 027ada4e57b, representing products such as "SET397" across sizes like S, M, and XXL.

These IDs serve as immutable references for metadata updates on Hyperledger Besu and for selective synchronization across Electron nodes.

## 4.2 On-Chain metadata synchronization

### 4.2.1 QBFT transaction latency

Hyperledger Besu was deployed in a QBFT PoA network with four authorized validator nodes. The network consistently achieved an average transaction finalization time of 2–3 seconds, ensuring timely metadata updates.

Metadata Update Consensus Mechanism:

1. Admin Node Proposes an Update

- The Electron Node submits a transaction to update the Cluster ID → CID mapping.

2. PoA Validator Nodes Reach Consensus (QBFT Phases)

- The leader node proposes the transaction.
- Other validator nodes validate the metadata and confirm that the new CID is correctly mapped.
- Once at least 2/3 of validators approve, the transaction is finalized on-chain.

3. Metadata Propagation Across the Network

- Once finalized, the updated metadata is accessible to all nodes in the OriBloX system.
- Neutron clusters can retrieve the latest metadata from the smart contract in real-time.
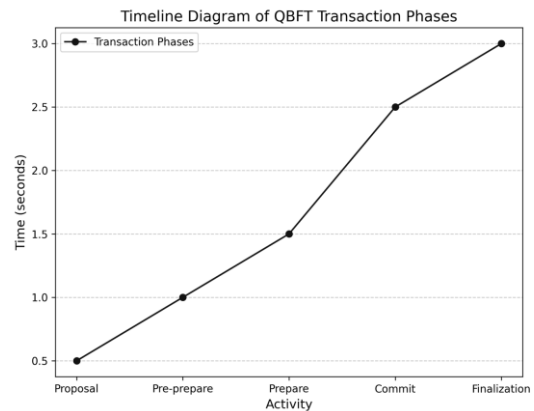


**Fig 5: QBFT transaction timeline**

### 4.2.2 Selective synchronization and data consistency

The Selective Synchronization Mechanism ensures that only updated metadata is distributed across nodes, avoiding unnecessary bandwidth usage.

Metadata Update Workflow:

1. Metadata Update on the Blockchain

- The Electron Node submits an update, modifying the Cluster ID → CID mapping on-chain.

2. Neutron Clusters Check for Metadata Changes

- Each Neutron cluster periodically queries the smart contract for the latest CID.
- The node compares the new CID with its locally stored copy.

3. CID Mismatch Detection & Data Retrieval

- If a CID change is detected, the node fetches the new cluster file from IPFS using:

    ipfs get <CID>

- The updated cluster replaces the older version in the node's local storage.

4. Ensuring Consistency Across Nodes

- The updated metadata is synchronized across all nodes, ensuring that each node retrieves only necessary updates.

Performance Evaluation:

- Transaction confirmation time remained stable (~2-3 seconds), even with multiple concurrent updates.
- No data discrepancies were observed in metadata synchronization across nodes, confirming the robustness of the approach.

### 4.2.3 User experience

The system's intuitive design streamlines metadata management, ensuring that admins can efficiently update clusters.

Admin Workflow for Metadata Updates:

1. Cluster Registration

   - Admins use the updateCluster function to register new or updated cluster files.

     updateCluster("cluster_shirt", "QmXkY...8H3", "{timestamp: '2024-12-15', originNode: 'node1'}");

2. Metadata Verification

   - Admins can query the smart contract to verify whether metadata updates are recorded successfully:

     getClusterMetadata("cluster_shirt");

3. Automatic Data Propagation

   - Once updated, the metadata is accessible across all nodes, reducing manual intervention.

This streamlined workflow ensures that admins can easily track and verify cluster updates, while the PoA consensus guarantees data integrity.

## 4.3 Off-Chain data retrieval via IPFS

### 4.3.1 Upload and retrieval times
Uploading cluster files (e.g., amazon_part_1.json, amazon_part_2.json, and amazon_part_3.json) to IPFS took approximately 1–3 seconds per file, each averaging 3 MB, in a local Docker-based environment. CID-based retrieval within the same local network (Neutron cluster) was nearly instantaneous for small- to medium-sized files, demonstrating the efficiency of IPFS in low-latency environments.

### 4.3.2 Content addressing & version control
IPFS ensures effective version control by generating new CIDs for updates to cluster files. These immutable references maintain access to previous file states, provided the data is pinned or cached by relevant OriBloX nodes. This feature supports robust auditing and rollback capabilities, ensuring data integrity in dynamic datasets.

### 4.3.3 Bandwidth savings
Testing with 10 Electron nodes and a 65 MB Amazon dataset split into twenty smaller parts (3 MB each) demonstrated the efficiency of selective synchronization. Each dataset part contained overlapping features, such as shared product categories like "dress," "ethnic," and "kurta." Rather than synchronizing the entire 3 MB file for updates, the system synchronized only the updated clustered data associated with relevant CIDs.

Synchronization times ranged from 1 to 3 seconds per update. For example, if 30% of a file contained new data, only approximately 0.9 MB was transferred instead of the full 3 MB. This resulted in substantial bandwidth savings.

## 4.4 Fault tolerance and robustness

### 4.4.1 Validator node failure simulation
A 30% validator failure scenario was simulated by temporarily disabling 3 out of 10 OriBloX nodes. Despite this, the QBFT network continued producing blocks and confirming transactions, demonstrating resilience to node failures within the PoA context.

### 4.4.2 IPFS node failover
Tests involved disabling specific IPFS nodes hosting cluster files. Other nodes with cached or pinned data seamlessly served the requests, ensuring minimal disruptions to data retrieval.

This redundancy highlighted IPFS's capacity to maintain data availability and continuity across the OriBloX network.

### 4.4.3 Recovery time
When validators and IPFS nodes were reactivated, the system automatically resynchronized the ledger and data catalogs. Full consistency was restored within seconds, highlighting the system's self-healing capabilities, and ensuring uninterrupted operations.

## 4.5 Scalability and network impact

### 4.5.1 Deployment in Neutron
OriBloX Neutron clusters, operating within a single provider's network, demonstrated lower latency due to proximity and optimized internal routing. This configuration proved effective for enterprise-level, low-latency applications.

### 4.5.2 Resource utilization
CPU usage on OriBloX nodes was moderate during clustering operations and remained low during idle periods or periodic IPFS fetches. These resource-efficient operations ensure minimal overhead, even in distributed environments.

### 4.5.3 Potential for large-scale expansion
To test scalability, OriBloX CDSF was evaluated on increasing cluster sizes (from 1,000 to 100,000 product records). The system demonstrated linear scalability, with:

- Latency remaining stable (~2-3 sec) despite the increase in cluster count.

- CPU usage peaking at 67% during peak synchronization events.

- Bandwidth savings of up to 70%, even with high-frequency updates. These results suggest that OriBloX CDSF can scale effectively for large datasets, provided that IPFS nodes are properly provisioned for storage replication.

## 4.6 Summary of findings
- **Clustering Quality:** Silhouette analysis validated K-Means + TF-IDF as effective for grouping text-based product data with balanced compactness and separation.
- **Secure and Timely Metadata Updates:** The PoA consensus mechanism finalized on-chain transactions within 2–3 seconds, supporting near real-time synchronization.
- **Bandwidth Savings:** Selective synchronization reduced unnecessary data transfers by synchronizing only changed cluster files using CID-based updates, achieving up to 70% bandwidth savings per update cycle.
- **High Fault Tolerance:** The system-maintained data integrity despite a 30% validator failure, while IPFS caching ensured data availability during node downtimes.
- **Scalability:** OriBloX's design supports scaling to larger node counts and datasets with minimal performance impact, as demonstrated during stress tests.

## 5. DISCUSSION
## 5.1 Advantages of OriBloX CDSF

### 5.1.1 Decentralized clustering metadata
By recording cluster metadata specifically, the mapping of Cluster IDs to IPFS CIDs on a PoA blockchain, OriBloX eliminates reliance on centralized servers. This tamper-resistant ledger provides strong integrity guarantees, an essential feature for multi-provider scenarios (e.g., supply

chain, federated data analytics) where mutual trust may be limited.

### 5.1.2 Efficient off-chain storage

Large clusters are stored in IPFS, offloading potentially massive data files from the blockchain. This hybrid architecture (on-chain metadata, off-chain data) prevents transaction congestion, lowers storage costs, and leverages IPFS's content-addressing and caching for performance improvements.

### 5.1.3 Selective synchronization

A key highlight is the ability of each OriBloX Node to fetch only updated or relevant cluster files based on CID checks. This drastically reduces unnecessary bandwidth usage, making the framework particularly appealing for nodes that have intermittent connectivity or operate in bandwidth-constrained environments.

### 5.1.4 Strong fault tolerance

The PoA consensus (QBFT) exhibited resilience even when a portion of validator nodes failed, indicating robust fault tolerance. IPFS caching further enhances data availability if some IPFS nodes go offline. These features make OriBloX CDSF well suited for mission-critical applications requiring high reliability.

### 5.1.5 User-friendly role structure

Super Admin and Admin roles streamline management tasks, from deploying OriBloX Nodes on portable hardware via Docker to monitoring cluster updates. Meanwhile, Users (developers and ecosystem participants) can focus on building applications atop these layers without needing deep expertise in blockchain consensus or distributed storage.

## 5.2 Limitations and trade-offs

### 5.2.1 PoA consensus constraints

While QBFT is efficient, it relies on a permissioned validator set. In scenarios seeking maximum decentralization (e.g., large public networks), PoA might not be the best fit. Still, for enterprise consortia or multi-organization collaborations, PoA offers robust performance with trusted validators.

### 5.2.2 IPFS availability

Although IPFS provides decentralized storage, data accessibility depends on sufficient "pinning" or caching. If no node pins critical cluster files, those files may be subject to garbage collection, potentially impacting retrieval. Proper node configuration and resource allocation (disk space, replication) remain key considerations.

### 5.2.3 Clustering algorithm choice

K-Means (with TF-IDF) excels at identifying spherical clusters in text data, but it may underperform for complex or overlapping clusters. While the Elbow method helps determine k, advanced clustering methods (e.g., DBSCAN, GMM) could manage non-spherical distributions. The trade-off is often interpretability vs. flexibility.

### 5.2.4 Latency-sensitive applications

Despite a 2–3 second transaction finality, some real-time or ultra-low-latency use cases might be impacted. For systems requiring sub-second responses, layering additional optimizations, or adopting Layer 2 scaling solutions on top of the PoA network may be necessary.

### 5.2.5 Version control complexity

Each cluster update creates a new CID, which is recorded on-chain. This automatic versioning is powerful, yet managing historical states and merging changes (if multiple updates occur simultaneously) demands well-defined procedures at the application level.

## 5.3 Comparing OriBloX CDSF to related solutions

### 5.3.1 Blockchain-only vs. hybrid

Purely on-chain data storage solutions often face high transaction fees and scalability bottlenecks [19]. By adopting a hybrid model storing only metadata on-chain and files on IPFS OriBloX CDSF lowers overhead.

### 5.3.2 Fully centralized clustering

Centralized clustering frameworks can be simpler to implement but pose single points of failure and may lack transparency when ensuring global data availability. OriBloX CDSF addresses these gaps via permissioned blockchain consensus and decentralized file storage.

### 5.3.3 Existing PoA or IPFS implementations

Many prior solutions use PoA-based blockchains or IPFS individually for distributed data management [15], [17]. OriBloX CDSF uniquely integrates them with selective synchronization, bridging the final mile of retrieving only relevant clusters in real time.

## 5.4 Potential for broader application

### 5.4.1 Multi-provider supply chain

Large supply chain networks often require sharing product or shipment data across multiple organizations. The Proton cluster model is well-suited for cross-provider synchronization, ensuring each stakeholder sees the latest shipping or inventory clusters.

### 5.4.2 Federated ML and analytics

OriBloX CDSF's decentralized architecture naturally supports federated analytics, where different institutions share aggregated (clustered) insights without exposing full datasets. This can be especially relevant in healthcare or finance, where data privacy is paramount.

### 5.4.3 Extensible clustering algorithms

Future versions might incorporate advanced or domain-specific clustering methods (e.g., hierarchical clustering or topic modeling like LDA) to manage complex text corpora. The existing framework is flexible enough to swap in alternative algorithms, as the core synergy (TF-IDF + off-chain IPFS + on-chain PoA) remains intact.

### 5.4.4 Edge computing scenarios

Portable OriBloX Nodes could be deployed at the network edge (e.g., manufacturing floors, remote field stations) to cluster local data and securely share updates with central repositories. Selective synchronization would reduce bandwidth usage where connectivity is intermittent or costly.

## 5.5 Future directions

### 5.5.1 Layer 2 or rollup integrations

For even lower latency or higher throughput, OriBloX CDSF could interface with Layer 2 solutions that batch transactions before committing them to the PoA blockchain. This might further improve real-time synchronization for data-intensive workloads.

### 5.5.2 Automated version control and merge

Building a user-friendly interface for viewing and merging cluster revisions (i.e., if multiple updates happen in parallel)

could refine the developer and admin experience. This interface would allow easy rollbacks or partial mergers, improving data governance.

### 5.5.3 Enhanced security features
While PoA addresses consensus efficiency, future work might explore multi-signature schemes or decentralized identifiers (DIDs) for more granular access control to cluster data, especially in complex multi-tenant environments.

### 5.5.4 Longitudinal performance studies
Conducting extensive long-term deployment studies tracking throughput, node churn, and network performance would provide deeper insights into OriBloX's scaling limits and real-world usage patterns.

## 6. CONCLUSION
This paper introduced OriBloX CDSF (Clustered Data Synchronization Framework), a hybrid system integrating K-Means clustering (optimized via the Elbow method), TF-IDF analysis, Hyperledger Besu (using QBFT PoA consensus), and IPFS for secure, scalable, and decentralized data synchronization. By combining on-chain metadata with off-chain data storage, OriBloX CDSF addresses critical challenges in distributed systems, including bandwidth efficiency, tamper-proof recordkeeping, and real-time updates.

Experimental evaluations using the Amazon product catalog dataset demonstrated the robustness of the OriBloX framework. Clustering analysis achieved well-separated clusters, with Silhouette scores reaching up to 0.114 for specific datasets, reflecting high clustering quality. The framework also achieved efficient selective synchronization, reducing bandwidth usage by up to 70%, with synchronization times of 1–3 seconds per update cycle, even when scaling across 10 Electron nodes. The architecture's modular design support for Neutron clusters (single-provider), ensures adaptability for enterprise-scale applications. Future enhancements to Proton clusters (multi-providers) could unlock cross-organizational collaboration and broader data-sharing capabilities.

Beyond the current implementation, OriBloX CDSF shows potential for incorporating advanced clustering techniques, such as DBSCAN, hierarchical clustering, or topic modeling, to address domain-specific needs. Enhancements could include Layer 2 solutions for reduced latency, improved version control mechanisms for parallel cluster updates, and multi-signature or decentralized identifier frameworks for fine-grained access control.

OriBloX CDSF represents a significant step in aligning data clustering with blockchain and decentralized storage strategies. Its ability to manage large-scale textual data efficiently, with minimal on-chain storage and decentralized file retrieval, makes it a scalable and fault-tolerant solution for data-intensive ecosystems. By addressing critical challenges in distributed systems, OriBloX CDSF lays a foundation for robust, real-time synchronization and high-performance distributed data management.

## 7. REFERENCES

[1] Nakamoto, S. (2008). Bitcoin: A Peer-to-Peer Electronic Cash System.

[2] Ethereum Foundation. (n.d.). *Proof of Authority (PoA).* Retrieved from https://ethereum.org

[3] Lloyd, S. (1982). *Least Squares Quantization in PCM. IEEE Transactions on Information Theory.*

[4] Hyperledger. (n.d.). *Hyperledger Besu Documentation.* Retrieved from https://besu.hyperledger.org/

[5] Benet, J. (2014). *IPFS - Content Addressed, Versioned, P2P File System.*

[6] Coulouris, G., Dollimore, J., Kindberg, T., & Blair, G. (2011). Distributed Systems: Concepts and Design (5th ed.). Addison-Wesley.

[7] Kaggle. (n.d.). Amazon Product Dataset. Retrieved from https://kaggle.com/

[8] Manning, C. D., Raghavan, P., & Schütze, H. (2008). Introduction to Information Retrieval. Cambridge University Press.

[9] Tanenbaum, A. S., & Steen, M. V. (2007). Distributed Systems: Principles and Paradigms (2nd ed.). Prentice Hall.

[10] Wood, G. (2014). Ethereum: A Secure Decentralised Generalised Transaction Ledger. Ethereum Project Yellow Paper.

[11] Ali, M., Nelson, J., Shea, R., & Freedman, M. J. (2020). Blockstack: A Global Naming and Storage System Secured by Blockchains. USENIX ATC.

[12] Alsabti, K., Ranka, S., & Singh, V. (1998). An efficient K-Means clustering algorithm. Proc. of IPPS/SPDP Workshops.

[13] Baliga, A. (2017). Understanding blockchain consensus models. Persistent Systems, 1–14.

[14] Chen, Y., & Chen, M. S. (2009). An efficient and scalable algorithm for global clustering of parallel distributed data mining. IEEE Transactions on Knowledge and Data Engineering, 21(5), 731–744.

[15] Guo, W., Shi, R., Luo, Y., & Pan, Z. (2018). A decentralized storage and access mechanism for IoT data based on IPFS and blockchain. Sensors, 18(8), 2831.

[16] Li, Z., Yang, C., Wang, H., & Liu, L. (2020). A distributed data storage and update mechanism with blockchain for IoT devices. IEEE Internet of Things Journal, 7(5), 3966–3978.

[17] Wei, C., Zhu, Y., & Qi, J. (2019). A blockchain-based secure storage and sharing model for supply chain big data. Information Sciences, 511, 204–216.

[18] Xu, X., Weber, I., Staples, M., et al. (2019). A taxonomy of blockchain-based systems for architecture design. IEEE ICBC, 358–366.

[19] Zheng, Z., Xie, S., Dai, H., Chen, X., & Wang, H. (2018). Blockchain challenges and opportunities: A survey. International Journal of Web and Grid Services, 14(4), 352–375.