# Analysis of Indexing Structures for Immutable Data

Cong Yue[†], Zhongle Xie[†], Meihui Zhang[§], Gang Chen[‡],
Beng Chin Ooi[†], Sheng Wang[⋆], Xiaokui Xiao[†]

[†] National University of Singapore
[§] Beijing Institute of Technology
[‡] Zhejiang University
[⋆] Alibaba Group

## ABSTRACT

In emerging applications such as blockchains and collaborative data analytics, there are strong demands for data immutability, multi-version accesses, and tamper-evident controls. To provide efficient support for lookup and merge operations, three new index structures for immutable data, namely Merkle Patricia Trie (MPT), Merkle Bucket Tree (MBT), and Pattern-Oriented-Split Tree (POS-Tree), have been proposed. Although these structures have been adopted in real applications, there is no systematic evaluation of their pros and cons in the literature, making it difficult for practitioners to choose the right index structure for their applications.

To alleviate the above problem, we present a comprehensive analysis of the existing index structures for immutable data, and evaluate both their asymptotic and empirical performance. Specifically, we show that MPT, MBT, and POS-Tree are all instances of a recently proposed framework, dubbed *Structurally Invariant and Reusable Indexes (SIRI)*. We propose to evaluate the SIRI instances on their index performance and deduplication capability. We establish the worst-case guarantees of each index, and experimentally evaluate all indexes in a wide variety of settings. Based on our theoretical and empirical analysis, we conclude that POS-Tree is a favorable choice for indexing immutable data.

## 1 INTRODUCTION

Accurate history of data is required for auditing and tracking purposes in numerous practice settings. In addition, data in the cloud is often vulnerable to malicious tampering. To support data lineage verification and mitigate malicious data
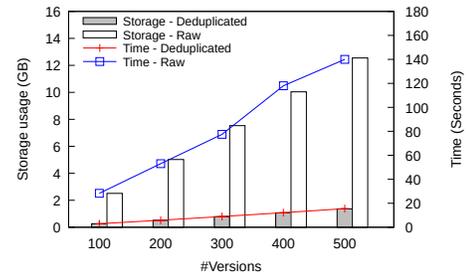
**Figure 1: Data storage and transmission time improved by deduplication**

manipulation, data immutability is essential for applications, such as banking transactions and emerging decentralized applications (e.g., blockchain, digital banking, and collaborative analytics). From the data management perspective, data immutability leads to two major challenges.

First, it is challenging to cope with the ever-increasing volume of data caused by immutability. An example is the sharing and storage of the data for healthcare analytics. Data scientists and clinicians often make relevant copies of current and historical data in the process of data analysis, cleansing, and curation. Such replicated copies could consume an enormous amount of space and network resources. To illustrate, let us consider a dataset that has 100,000 records initially, and it receives 1,000 record updates in each modification. Figure 1 shows the space and time required to handle the increasing number of versions[1]. Observe that (i) the space and time overheads are significant if all versions are stored separately, and (ii) such overheads could be considerably reduced if we can *deduplicate* the records in different versions.

The second challenge is that in case a piece of data is tampered with (e.g., malicious manipulation of crypto-currency wallets or unauthorized modifications of patients' lab test data), we have to detect it promptly. To address this challenge, the system needs to incorporate tamper-resistant techniques to support the authentication and recovery of data, to ensure data immutability. Towards this end, a typical

---

[1]Run with Intel(R) Xeon(R) E5-1620 v3 CPU and 1 Gigabit Ethernet card.

approach is to adopt cryptographic methods for tamper mitigation, which, however, considerably complicates the system design.

Most existing data management solutions tackle the above two challenges separately, using independent orthogonal methods. In particular, they typically (i) ensure tamper evidence using cryptographic fingerprints and hash links [30], and (ii) achieve deduplication with delta encoding [23, 27]. Such decoupled design incurs unnecessary overheads that could severely degrade the system performance. For example, in state-of-the-art blockchain systems such as Ethereum [3] and Hyperledger [7], tamper evidence is externally defined and computed on top of the underlying key-value store (e.g., LevelDB [8] or RocksDB [11]), which leads to considerable processing costs. In addition, delta-encoding-based deduplication (e.g., in Decibel [27]) requires a reconstruction phase before an object can be accessed, which renders data accessing rather inefficient.

Motivated by the above issues, recent work [7, 44, 45] has explored data management methods to provide native supports for both tamper evidence and deduplication features. This results in three new index structures for immutable data, namely, *Merkle Patricia Trie (MPT)* [45], *Merkle Bucket Tree (MBT)* [7], and *Pattern-Oriented-Split Tree (POS-Tree)* [44]. To the best of our knowledge, there is no systematic comparison of these three index structures in the literature, and the characteristics of each structure are not fully understood. This renders it difficult for practitioners to choose the right index structure for their applications.

To fill the aforementioned gap, this paper presents a comprehensive analysis of MPT, MBT, and POS-Tree. Specifically, we make the following contributions:

- We show that MPT, MBT, and POS-Tree are all instances of a recently proposed framework, named Structurally Invariant and Reusable Indexes (*SIRI*) [44]. Based on this, we identify the common characteristics of them in terms of tamper evidence and deduplication.
- We propose a benchmarking scheme to evaluate *SIRI* instances based on five essential metrics: their efficiency for four index operations (i.e., lookup, update, comparison, and merge), as well as their *deduplication ratios*, which is a new metric that we formulate to quantify each index's deduplication effectiveness. We establish the worst-case guarantee of each index in terms of these five metrics.
- We experimentally evaluate all three indexes in a variety of settings. We demonstrate that they perform much better than conventional indexes in terms of the effectiveness of deduplication. Based on our experimental results, we conclude that POS-Tree is a favorable choice for indexing immutable data.

The rest of the paper is organized as follows. Section 2 presents the background and prior researches on the core application properties. Section 3 presents *SIRI*, along with an extended discussion on its significance and the explanation of three *SIRI* representatives. A theoretical analysis is conducted in Section 4 to reveal the operational bounds of *SIRI* while the experimental evaluation is reported in Section 5. We conclude this paper in Section 6.

## 2 RELATED WORK

We first discuss the background and several primary motivations leading to the definition of *SIRI*.

### 2.1 Versioning and Immutability

Data versioning has been widely employed for tolerating failures, errors, and intrusions, as well as for analysis of data modification history. ElephantFS [38] is one of the first-generation file systems with built-in multi-version support. Successor systems like S4 [42], CVFS [40], RepareStore [47] and OceanStore [25] improve the early design by maintaining all versions in full scope and upon each update operation. In databases, data versioning techniques are used for transactional data access. Postgres [17, 41], for example, achieved comparable performance to the database systems without versioning support. Fastrek [19] enhanced Postgres with intrusion tolerance by maintaining an inter-transaction dependency graph based on the versioned data and relying on the graph to resolve data access conflicts. Multi-versioning is also used to provide snapshot isolation in database systems [16, 33] although such systems usually do not store the full history versions. To directly access multi-versioned data, a number of multi-version data structures can be applied from the literature, such as multi-version B-tree [26, 35], temporal hashing [24] and persistent data structures [31, 34].

Immutable data are becoming versatile in emerging applications. For example, blockchains [3, 7, 20, 21, 30, 36] maintain immutable ledgers, which keep all historical versions of the system status. Similarly, collaborative applications [6, 18] maintain the whole evolutionary history of datasets and derived analytic results, which enables provenance-related functionalities, such as tracking, branching, and rollback. A direct consequence of data immutability is that all stored data are inherently multi-versioned upon being amended. There exist a wide range of storage systems handling such data in either a linear manner, such as multi-version file systems [38, 40, 42] and temporal databases [14, 37, 43], or a non-linear manner, such as version control systems including git [4], svn [12] and mercurial [9], and collaborative management databases including Decibel [27] and OrpheusDB [23]. Git and Git-like systems are also used to manage the history and branches of datasets to achieve efficient query and space utilization [5, 18, 39].

## 2.2 Data-Level Deduplication

Deduplication approaches have been proposed to reduce the overhead of storage consumption when maintaining multi-versioned data. For example, Decibel [27] uses *delta encoding*, whereby the system only stores the differences, called delta, between the new version and the previous version of data. Consequently, it is effective to manage data versions when the deltas are small, despite the extra cost incurred during data retrieval for reconstructing the specified version of data. However, it is ineffective in removing duplicates among non-consecutive versions or different branches of the data. Though some algorithms choose a more precedent version that has the smallest differences as the parent to improve the efficiency of the deduplication, it involves additional complexity to reconstruct a version.

To enable the removal of duplicates among any data versions, *chunk-based deduplication* can be applied. Unlike delta encoding, this approach works across independent objects. It is widely used in file systems [32, 46], and is a core principle of git. In this approach, files are divided into *chunks*, each of which is given a unique identifier calculated from algorithms like collision-resistant hashing. *chunks* with the same identifier can be eliminated. Chunk-based deduplication is highly effective in removing duplicates for large files that are rarely modified. In case an update leads to a change of all subsequent chunks, i.e., the boundary-shifting problem [22], content-defined chunking [29] can be leveraged to avoid expensive re-chunking.

## 2.3 Tamper Evidence

Applications such as digital banking [13] and blockchain [3, 7, 30] demand the system should maintain the accurate history of data, protect their data from malicious tampering, and trigger alerts when malicious tampering occur. To serve such purposes, verifiable databases (i.e., Concerto [15], QLDB [1]) and blockchain services (i.e., Microsoft Azure Blockchain [2]) often use cryptographic hash functions (e.g., SHA) and Merkle trees [28] to verify the data integrity. *SIRI*, with the built-in support for tamper evidence, is a good candidate for the above systems.

A Merkle tree is a tree of hashes, where the leaf nodes are the cryptographic hashes calculated from blocks of data while the non-leaf nodes are the hashes of their immediate children. The root hash is also called the "digest" of the data. To verify a record, it requires a "proof" of data, which contains the nodes on the path to the root. The new root hash is recalculated recursively and equality is checked with the previously saved digest.
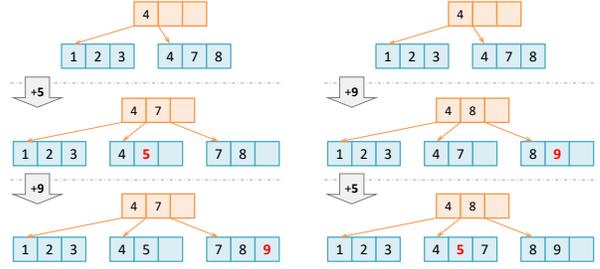


Figure 2: Two B$^+$-trees containing the same entries but with different internal structures [44]

# 3 STRUCTURALLY INVARIANT AND REUSABLE INDEXES

Structurally Invariant and Reusable Indexes (*SIRI*) are a new family of indexes recently proposed [44] to efficiently support tamper evidence and effective deduplication.

## 3.1 Background and Notations

In addition to basic lookup and update operations, the ultimate goal of *SIRI* is to provide native data versioning, deduplication and tamper evidence features. Consequently, data pages in *SIRI* must not only support efficient deduplication (to tackle the amount of replication arising from versioning) but also cryptographic hashing (to facilitate tamper evidence).

To better elaborate the *SIRI* candidates, we use the following notations in the remaining of this paper. The indexing dataset is denoted as $\mathcal{D} = \{D_0, D_1, ..., D_n\}$ where $D_i$ represents its $i$-th version. $\mathcal{I}$ is employed to represent *SIRI* structures, and $I$ stands for one of its instances. The key set stored in $I$ is set as $R(I) = \{r_1, r_2, ..., r_n\}$, where $r_i$ denotes the $i$-th key. $P(I) = \{p_1, p_2, ..., p_n\}$ stands for the internal node set of $I$, where $p_i$ represents the $i$-th node.

## 3.2 Formal Definition

We provide a formal and precise definition of *SIRI* adapted from [44] as follows.

*Definition 3.1.* An index class $\mathcal{I}$ belongs to *SIRI* if it has the following properties:

(1) *Structurally Invariant.* If $I$ and $I'$ are two instances of $\mathcal{I}$, then $P(I) = P(I') \iff R(I) = R(I')$.
(2) *Recursively Identical.* If $I$ and $I'$ are two instances of $\mathcal{I}$ and $R(I) = R(I') + r$, where $r \notin R(I')$, then $|P(I) \cap P(I')| \gg |P(I) - P(I')|$.
(3) *Universally Reusable.* For any instance $I$ of $\mathcal{I}$, there always exists node $p \in P(I)$ and another instance $I'$ such that $|P(I')| > |P(I)|)$ and $p \in P(I')$.

Definition 3.1 states that *SIRI* must possess three properties. The first property, *Structurally Invariant*, ensures that the order of update operations does not affect the internal

structure of the index, while the second property, *Recursively Identical*, guarantees the efficiency when constructing a large instance from small ones. The third property, *Universally Reusable*, secures that the nodes of the index could be shared among different instances. In practice, these properties can be exploited to make *SIRI* time- and space-efficient.

## 3.3 Extended Discussion

*Recursively Identical* and *Universally Reusable* are both aimed at making the pages share-able among various instances. However, they focus on different aspects. The former attribute concentrates on providing performance improvement when designing the indexes – updates do not bring in harmful impacts since the performance is often dominated by accessing a vast number of shared pages. The latter is to secure the theoretical boundary of *SIRI*'s performance. The higher the ratio of shared pages each instance gets, the better performance *SIRI* could reach in terms of deduplication. In the limiting case, where the dataset and indexing operations are infinite, every page in a *SIRI* instance could find its copy used by other instances.

It is non-trivial to construct a *SIRI* instance from conventional structures. Take the multi-way search tree as an example. Such a structure is *Recursively Identical* since only a small part of nodes is changed in the new version of the instance when an update operation is performed. Further, the usage of copy-on-write implementation naturally enables node sharing among versions and branches. Hence, it can be *Globally Reusable* when applying this technique. However, it may not be *Structurally Invariant*. Take B$^+$-tree as an example, Figure 2 illustrates that identical sets of items may lead to variant structures. Meanwhile, hash tables are not *Recursively Identical* when they require periodical reconstructions as the entire structure may be updated and none of the nodes can be reused.

Surprisingly, tries, or radix trees, can meet all the three properties with copy-on-write implementation. Firstly, they are *Structurally Invariant* since the position of the node only depends on the sequence of the stored key bytes and consequently, the same set of keys always leads to the same tree structure. Secondly, being a multi-way search tree, they can be *Recursively Identical* and *Globally Reusable* as mentioned above. However, they may end up in higher tree heights, leading to poor performance caused by increasing traversal cost, as shown in Section 5.

Due to the appearance of the three properties, the adoption of the aforementioned data-level deduplication approaches can be seamlessly applied in index-level for *SIRI*. The identical pages from different index instances for multiple versions of the data can be shared and therefore, the system can persist only one copy to save space. Another benefit of applying index-level deduplication is that the system can access a
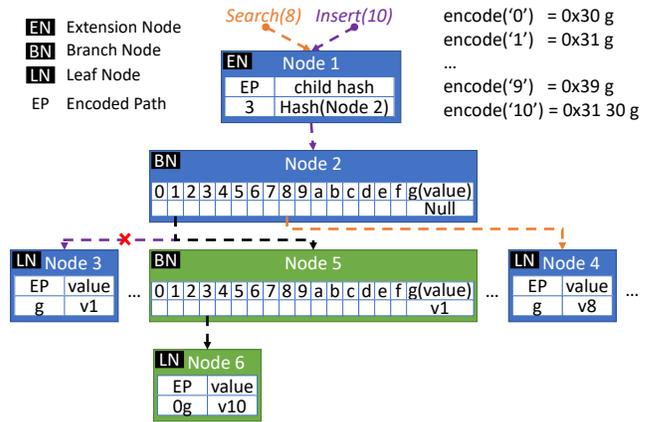


**Figure 3: Merkle Patricia Trie (MPT)**

version of the data directly from different indexes and data pages instead of experiencing a reconstruction phase from the deltas. Overall, the nature of *SIRI* renders effective detection and removal of duplicates without prohibitive efforts, which conventional index structures can hardly offer.

## 3.4 SIRI Representatives

In this section, we elaborate on the three representatives of *SIRI*, namely MPT, MBT, and POS-Tree. As mentioned in section 3.3, all representatives are *Recursively Identical* and *Globally Reusable* being multi-way search trees and leveraging a copy-on-write implementation of their nodes. Meanwhile, they are *Structurally Invariant* as stated in Section 3.4.1, 3.4.2, and 3.4.3.

### 3.4.1 Merkle Patricia Trie.

Merkle Patricia Trie (MPT) is a radix tree with cryptographic authentication. Similar to the traditional radix tree, the key is split into sequential characters, namely nibbles. There are four types of nodes in MPT, namely *branch*, *leaf*, *extension* and *null*. The structures of those nodes are illustrated in Figure 3: (1) *branch* node consists of a 16-element array and a value. Each element, called "branch", of the array is indexing a corresponding child node and stores a nibble. (2) *leaf* node contains a byte string, i.e., a compressed path called "encodedPath", and a value. (3) *extension* node also contains encodedPath and a pointer to the next node. (4) *null* node includes an empty string indicating that the node contains nothing. Similar to Merkle Tree, the whole MPT can be rolled up to a single cryptographic hash for tamper evidence. The most well-known usage of this data structure is in Ethereum [3], one of the largest blockchain systems in the world.

**Lookup.** The lookup procedure for key "8" is illustrated in Figure 3. The key is first encoded as "0x38 g". Then, each character of the encoded key is used to match with the encodedPath in an extension node, or to select the path in a branch node, from left to right. For this example, the first
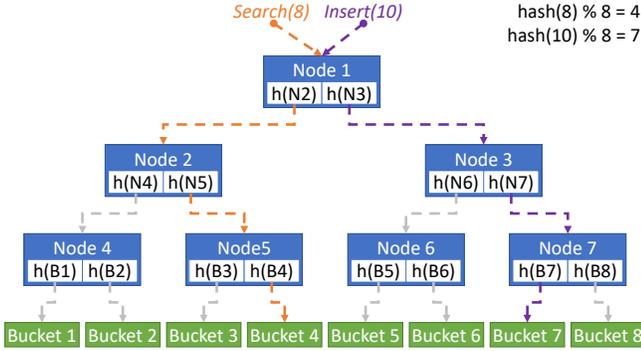
**Figure 4: Merkle Bucket Tree (MBT)**

character "3" matches the root node's encodedPath, therefore, it navigates to its child, Node 2. Then it takes the branch "8" since "8" equals to the second character in the encoded key. Finally, the traversal reaches the leaf node and ends with the value "v8" output.

**Insert.** To insert data in MPT, the index first locates the position of the given key as in the lookup operation. Once it reaches a *null* node, a leaf node containing the remaining part of the encoded key and the value is created. For example, in Figure 3, when we insert key "1" ("0x31 g"), if branch "1" in Node 2 is empty, a new node ("g", v1) is created and pointed by branch "1". In case there is a partial match at extension node, a new branch node at diverging byte is created, appended with original and new child. The insertion of key "10" in the figure can illustrate this procedure, where the path is diverged at Node 3. Hence, Node 3 is replaced by Node 5 with a newly created Node 6 attached.

### 3.4.2 Merkle Bucket Tree.

Merkle Bucket Tree (MBT) is a Merkle tree built on top of a hash table as shown in Figure 4. The bottom most level of MBT is a set of buckets and the cardinality of the bucket set is called *capacity*. Data entries are hashed to these buckets, and the entries within each bucket are arranged in sorted order. The internal nodes are formed by the cryptographic hashes computed from their intermediate children. The number of children an internal node has is called *fanout*. In MBT, *capacity* and *fanout* are pre-defined and cannot be changed in its life cycle.

**Lookup.** To perform an MBT index lookup, we first calculate the hash of the target key and obtain the index of the bucket where the data resides. Due to the copy-on-write restrictions, we are unable to fetch the bucket directly and hence we then use the bucket number to calculate the traversal path from the root node to the leaf node. The calculation is generally a trivial reverse simulation of the complete multiway search tree search algorithm. For example in Figure 4, key "8" falls into Bucket 4 after the hashing, and we accordingly get all node index on the path starting from the leaf.

Finally, we follow the path to reach the bucket. The records in the bucket are scanned using binary search to find the target key after the retrieval of the bucket node.

**Insert.** The insert operation of MBT undergoes similar procedures. It first performs a lookup to check the existence of the target key. For example, the inserting key "10" falls into Bucket 7. Then Bucket 7 is fetched following the lookup process, and the key is inserted to Bucket 7 in ascending order. Finally, the hashes of the bucket and the nodes are recalculated recursively.

The design of MBT undoubtedly takes the advantages of Merkle tree and the hash table. On the one hand, MBT offers tamper evidence with a low update cost since only the set of nodes lying on the lookup path needs to be recalculated. On the other hand, the data entries can be evenly distributed due to the nature of the hash buckets in the bottom level.

### 3.4.3 Pattern-Oriented-Split Tree.

Pattern-Oriented-Split Tree (POS-Tree) is a probabilistically balanced search tree proposed in [44]. The structure can be treated as a customized Merkle tree built upon pattern-aware partitions of the dataset, as shown in Figure 5. The bottom most data layer is an ordered sequence of data records. The records are partitioned into blocks using a sliding-window approach and such blocks form the leaf nodes. That is, for a byte sequence within a fixed-sized window, starting from the first byte of the data, a Rabin fingerprint is computed to match a certain boundary pattern. An example pattern can be the last 8 bits of Rabin fingerprint equaling to "1". The window shifts forward to repeat the process until it finds a match, where the node boundary is set to create the leaf node. The internal layers are formed by a sequence of split keys and cryptographic hashes of the nodes in the lower layer. Since the contents in the internal layers already contain hash values, we directly use the hashes to match the boundary pattern instead of repeatedly computing the hashes within a sliding window. Such strategy improves the performance of POS-Tree by reducing the number of hash computations, while preserving the randomness of chunking.

**Lookup.** The lookup procedure of POS-Tree is similar to B$^+$-tree. Starting from the root node, it performs binary search to locate the child node containing the target key. When it reaches the leaf node, a binary search is performed to find the exact key. As the example shown in Figure 5, the key "8" is fetched through the orange path. It goes through Node 2, which has a key range of (-∞, 351], and Node 4, which has a key range of (-∞, 89].

**Insert.** To perform an insert operation, POS-Tree first finds the position of the inserting key and then inserts it into the corresponding leaf node. Next, it starts the boundary detection from the first byte of the leaf node, and stops when
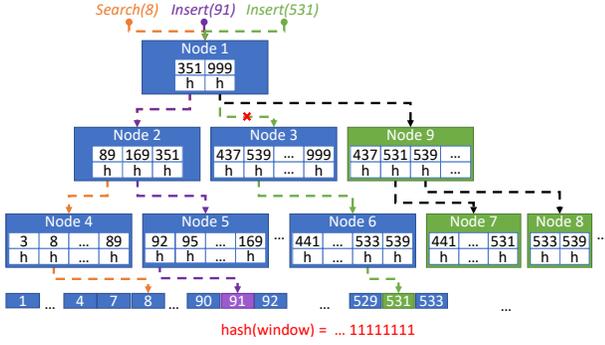
**Figure 5: Pattern-Oriented-Splitting Tree (POS-Tree)**
**Table 1: Notation table**

| Symbol | Description |
|---|---|
| $N$ | The total number of records |
| $m$ | The fanout of POS-Tree and MBT |
| $B$ | The capacity of MBT (Number of buckets) |
| $L$ | The key length of a record |
| $\delta$ | The number of different records between two versions |
| $\alpha$ | The ratio of the number of different records and the number of total records |
| $r$ | The average size of a record |
| $c$ | The size of cryptographic hash value |

detecting an existing boundary or reaching the last byte of the layer. For example, when insert key "91" into the tree shown in Figure 5, a boundary detection is performed from Node 5. It ends upon reaching the existing boundary of Node 5. Another instance in the figure is the insertion of key "531". A new boundary is found at element 531, and the traverse stops when finding the existing boundary of Node 6. Therefore, Node 6 splits into Node 7 and Node 8, and the new split keys are propagated to the parent node.

The pattern-aware partitioning of POS-Tree enhances the deduplication capabilities, and making the structure of the tree depending only on the data held. Such *Structurally Invariant* property supports efficient *diff* and *merge*. Moreover, the B$^+$-tree-like node structure enables efficient indexing by comparing the split keys to navigate the paths.

# 4 THEORETICAL ANALYSIS

In this section, we provide a comprehensive theoretical analysis of the three *SIRI* representatives discussed previously. We first calculate the theoretical bounds for common index operations like lookup and update, as well as complex operations needed by emerging applications including diff and merge. In addition, we define the *deduplication ratio* as a metric for measuring the efficiency of deduplication provided by *SIRI*s.

In the following subsections, $N$ denotes the maximum capacity of distinct keys in an index; $L$ denotes the maximum length of a key string; $B$ denotes the number of buckets in MBT. $m$ denotes the expected number of entries per page/node in MBT or POS-Tree. $\delta$ denotes the different records between two instances. $r$ denotes the average storage size for a sole record.

## 4.1 Operation Bounds

In this section, the bounds of common operations are calculated accordingly.

### 4.1.1 Index Lookup.

We first evaluate the lookup complexity of the three candidates.

- **MPT** – The traversal upon MPT is the same as a normal radix tree with path compaction. The computing bound for lookup is the maximum between $O(L)$ and $O(\log_m N)$. Since $L$ is often larger than $\log_m N$ in the real systems ($L$ equals to 64-byte in Ethereum's zeroith variable), the lookup complexity in MPT is $O(L)$ in most of the time.
- **MBT** – Unlike other structures having constant leaf node scanning time, the size of MBT's leaf node is $\frac{N}{B}$ and it therefore costs $O(\log_2 \frac{N}{B})$ with binary search to scan the node. As a result, the total lookup cost, consisting of node traversing part and leaf node scanning part, is $O(\log_m B + \log_2 \frac{N}{B})$.
- **POS-Tree** – Being a probabilistically balanced search tree, the complexity of POS-Tree is $O(\log_m N)$.

### 4.1.2 Index Update.

In all candidates, an update operation firstly incurs a lookup for the updating keys, and then leads to the creation of the copy of affected nodes and the calculation of their hash values. In our analysis, we treat the number of entries in each node $m$, the length of each record $r$ and the length of hashed value $len(h)$ as constant values. That is to say, the size of internal nodes $m \cdot len(h)$, and the size of leaf nodes $m \cdot r$ are constant unless explicitly stated. Therefore, the cost of new node creation and the crypto-hash function is constant, and the cost of the update mainly depends on the cost of the lookup.

The calculation and comparison among the three candidates are listed below:

- **MPT** – $\max(O(L), O(\log_m N))$. In most cases, the complexity of the update in MPT is $O(L)$.
- **MBT** – For leaf nodes, as the size increases linearly with $N$, the complexity of hash function and node copying is $O(\frac{N}{B})$. Hence, the complexity of the update in MBT is $O(\log_m B + \frac{N}{B})$.

- **POS-Tree** – Similar to crypto-hash function, the cost of the rolling hash function p for detecting node boundary is also constant. This results in the update complexity of $O(\log_m N)$.

### 4.1.3 Indexes Diff.

*Diff* is the operation that compares two index instances. It returns all records that are either present in only one index or different in both indexes. Therefore, *Diff* can be seen as multiple lookups in a naive implementation of the three candidates. The following bounds are calculated under this assumption. We directly give the results due to its triviality.

- **MPT** – $O(\delta \cdot L)$ or $O(\delta \cdot \log_m N)$. As discussed previously, in most cases the complexity is the former.
- **MBT** – $O(\delta \cdot (\log_m B + \frac{N}{B}))$
- **POS-Tree** – $O(\delta \cdot \log_m N)$

### 4.1.4 Indexes Merge.

*Merge* is the operation that combines all records from either indexes. The entire process of *Merge* contains two steps. The first step is to do a *Diff* operation between the instance to merge and the original instance, mark the different pages/nodes. The second step is to merge all the different nodes into the original instance. If there exist conflicts, namely a key in both instances with different values, the process must be interrupted and a selection strategy must be given by the end user to continue. The following calculation is based on the worst case when the merge process can be finished without interruption. Since the second step of the merge process is treated as $O(1)$ operations in our analysis, the complexity of the merge is dominated by the "diff" operation in the first step.

- **MPT** – $O(\delta \cdot L)$ or $O(\delta \cdot \log_m N)$. In most cases, the complexity should be $O(\delta \cdot L)$.
- **MBT** – $O(\delta \cdot (\log_m B + \frac{N}{B}))$
- **POS-Tree** – $O(\delta \cdot \log_m N)$

In the worst case, MPT has higher tree height than a balanced search tree, i.e., L ¿ $O(\log_m N)$, and therefore performs worse than POS-Tree. For MBT, the traverse cost $\log_m B$ is lower than other structures when in assumption of $B < N$ while the node scanning time $\log_2 \frac{N}{B}$ and creation time $\frac{N}{B}$ are dominating when $N >> B$. We can conclude from the table that POS-Tree is efficient in general cases, while MBT is a good choice when the dataset maintains a proper $\frac{N}{B}$ ratio.

## 4.2 Deduplication Ratio

Persistent (or immutable) data structures demand a large amount of space for maintaining all historical versions of data. To alleviate space consumption pressure, the feasibility of detecting and removing duplicated data portions plays a critical role. In this section, we aim to quantify the effectiveness of such properties in indexes by defining a measurement called *deduplication ratio*.

### 4.2.1 Definition.

Suppose there is a set of index instances $S = \{I_1, I_2, ...I_k\}$, and each $I_x$ is composed of a set of pages $P_x$. The byte size of a page $p$ is denoted as $byte(p)$, we can derive byte count of set $P$ as:

$$byte(P) = \sum_{p \in P} byte(p).$$

The *deduplication ratio $\eta$* of $S$ is defined as follows:

$$\eta(S) = 1 - \frac{byte(P_1 \cup P_2 \cup ... \cup P_k)}{byte(P_1) + byte(P_2) + ... + byte(P_k)},$$

or

$$\eta(S) = 1 - \frac{byte(\bigcup_{i=1}^{k} P_i)}{\sum_{j=1}^{k} byte(P_j)}.$$

The $\eta$ quantifies the effectiveness of page-level data deduplication (i.e., sharing) among related indexes. It is the ratio between the overall bytes that can be shared between different page sets and the total bytes used for all the page sets. With a high $\eta$, the storage is capable of managing massive "immutable" data versions without bearing space consumption pressure. In the following subsections, we will use this metric to evaluate the three candidates accordingly.

### 4.2.2 Continuous Differential Analysis.

In this part, we analyze a simple case that $S$ consists $n$ sequentially evolved indexes, i.e., the $i_{th}$ instance is derived from the $i-1_{th}$ instance. Each instance $S_i$ can be represented as a page set $P_i$ or a record set $R_i$. The analysis of more complicated scenarios is treated as our future work. To ease our analysis, we assume that each instance differs its predecessor by ratio $\alpha$ of a continuous key range $\delta$, such that:

$$|D_i| = \alpha \cdot |R_{i-1}|,$$

$$\forall k \in (R_i - D_i) \cup (R_{i-1} - D_i), (k < min(D_i) \vee k > max(D_i))$$

where $|X|$ denotes the record count in set $X$, and $min/max$ denotes the minimum/maximum key in a set.

In the following analysis, we consider two scenarios:

- Insertion of new records.

$$|R_i| = (1 + \alpha) \cdot |R_{i-1}|.$$

- Update of existing records.

$$|R_i| = |R_{i-1}|.$$

**Merkle Bucket Tree.** Since in MBT, the bucket size depends on the number of contained records, i.e.,

$$E = \frac{N}{B}.$$

We denote the number of affected nodes on level $x$ in MBT as $N_x$. Hence, the number of buckets (the leaf level) affected by $\alpha$ differential is expressed as:

$$N_0 = \alpha \cdot B.$$

We can roughly summarize the total number of affected tree nodes, $N_{tree}$, as following:

$$N_{tree} = \left\lceil \frac{B}{m} \right\rceil + \left\lceil \frac{B}{m^2} \right\rceil + \dots + \left\lceil \frac{B}{m^{\log_m B}} \right\rceil$$
$$\approx \frac{B}{m} + \frac{B}{m^2} + \dots + \frac{B}{m^{\log_m B}}$$
$$= \frac{B-1}{m-1}.$$

The number of affected nodes in the continuous update can be calculated as:

$$\sum_{i=1}^{\log_m B} N_i = \left\lceil \alpha \cdot \frac{B}{m} \right\rceil + \left\lceil \alpha \cdot \frac{B}{m^2} \right\rceil + \dots + \left\lceil \alpha \cdot \frac{B}{m^{\log_m B}} \right\rceil$$
$$\approx \alpha \cdot \frac{B}{m} + \alpha \cdot \frac{B}{m^2} + \dots + \alpha \cdot \frac{B}{m^{\log_m B}}$$
$$= \alpha \cdot \frac{B-1}{m-1}.$$

Thus, the deduplication ratio for MBT is:

$$\eta(MBT) = 1 - \frac{\text{bytes in modified nodes}}{\text{total bytes}}$$
$$= 1 -$$
$$\frac{\alpha \cdot N \cdot r + m \cdot c \cdot \sum_{i=1}^{\log_m B} N_i + N \cdot r + m \cdot c \cdot N_{tree}}{2 \cdot (N \cdot r + m \cdot c \cdot N_{tree})}$$
$$\approx \frac{1}{2} - \frac{\alpha \cdot N \cdot r + m \cdot c \cdot \alpha \cdot \frac{B-1}{m-1}}{2 \cdot (N \cdot r + m \cdot c \cdot \frac{B-1}{m-1})}$$
$$= \frac{1}{2} - \frac{\alpha}{2}.$$

Surprisingly, the deduplication ratio is highly related to $\alpha$ and has no direct connection with $B$ according to the analysis result.

**Merkle Patricia Trie.** In a simple MPT without any path compaction optimization, we have:

$$\eta(MPT) = 1 - \frac{\text{bytes in modified nodes}}{\text{total bytes}}$$
$$= \frac{1}{2} - \frac{\alpha \cdot N \cdot (L \cdot c + r)}{2 \cdot (N \cdot r + N \cdot \bar{L} \cdot c)},$$

which indicates that

$$\eta(MPT) \geq \frac{1}{2} - \frac{\alpha}{2} \ (L \geq \bar{L})$$

,

$$\eta(MPT) \leq \frac{1}{2} - \frac{\alpha}{2} \ (L \leq \bar{L})$$

**Table 2: Parameter table for experiments**

| Parameter | Value |
|---|---|
| Dataset size($10^4$) | 1, 2, 4, 8, <u>16</u>, 32, 64, 128, 256 |
| Batch size($10^3$) | 1, <u>2</u>, 4, 8, 16 |
| Overlap Ratio | <u>0</u>, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100 |
| Write Ratio(%) | <u>0</u>, 50, 100 |
| Zipfian parameter $\theta$ | <u>0</u>, 0.5, 0.9 |

Inferred from the result, $\eta(MPT)$ is affected by the distribution of stored keys since the length of the average length of the keys highly relate to the final deduplication ratio. In detail, the relationship between $L$ and $\bar{L}$ determines whether the deduplication ratio of MPT is greater than or less than that of MBT.

**POS-Tree.** Similar to MBT, the calculation is as follows:

$$\eta(POS) = 1 - \frac{\text{bytes in modified nodes}}{\text{total bytes}}$$
$$\approx \frac{1}{2} - \frac{\alpha \cdot N \cdot r + m \cdot c \cdot \sum_{j=1}^{\log_m \frac{N}{f}} N_j}{2 \cdot (N \cdot r + \frac{\frac{N}{f}-1}{m-1} \cdot m \cdot c)}$$
$$\approx \frac{1}{2} - \frac{\alpha \cdot N \cdot r + m \cdot c \cdot \alpha \cdot \frac{\frac{N}{f}-1}{m-1}}{2 \cdot (N \cdot r + \frac{\frac{N}{f}-1}{m-1} \cdot m \cdot c)}$$
$$= \frac{1}{2} - \frac{\alpha}{2}.$$

If we compare the analysis results of the three representatives, we can conclude that MPT has the best deduplication ratio under proper query workloads and datasets (meaning $L \geq \bar{L}$). Meanwhile, POS-Tree and MBT have equal bound for the deduplication ratio in this setting.

# 5 EXPERIMENTAL BENCHMARKING

In this section, we evaluate three *SIRI* representatives, namely POS-Tree, MBT and MPT, through different experiments. First, the throughput and the latency of the indexes are measured to have an overview of how these structures perform in general cases. Second, the storage consumption, the deduplication ratio and the node sharing ratio are evaluated to investigate the space efficiency among the candidates. Third, a breakdown analysis is given to show how each *SIRI* property affects the performance of the index. Finally, we integrate the structures in an existing database management system, Forkbase [44], to show how *SIRI* structures behave in real applications.

Our experiments are conducted on a server with Ubuntu 14.04, which is equipped with an Intel Xeon Processor E5-1650 processor (3.5GHz) and 32GB RAM. To fairly compare the efficiency of the index structures in terms of node quantity and size, we tune the size of each index node to be approximately 1 KB. For each experiment, the reported measurements are averaged over 5 runs.

## 5.1 Dataset

We use a synthesized YCSB dataset and two real-world datasets, Wikipedia data dump and Ethereum transaction data, to conduct a thorough evaluation of SIRI.

*5.1.1 YCSB.* We generate the key-value dataset using YCSB according to the parameters shown in Table 2. The lengths of the keys range from 5 bytes to 15 bytes, while the values have an average length of 256 bytes. The total number of records varies from 10,000 to 2,560,000. The dataset contains three types of workloads, read, write and mixed workload with 50% write operations. We use Zipfian distribution to simulate the scenarios where the data in the workload is skewed to different degrees, where the Zipfian parameter $\theta$ equals 0 represents all records have an equal possibility to be selected into the workload, while higher $\theta$ value means only a smaller range of the records have extremely high possibilities to be selected. We also generate the overlapped workloads to test the capability of deduplication with increasing similarity in the contents, as described in Section 5.4.2.

*5.1.2 WIKI.* The wiki dataset is real-world Wikipedia data dumps[2] of the extracted page abstracts. The key of the dataset is the URL of the Wikipedia page, the length of which ranges from 31 bytes to 298 bytes and has an average of 50 bytes. While the value of the dataset is the extracted abstract in plain text format, the length of which ranges from 1 byte to 1036 bytes, having an average of 96 bytes. We collect 6 data dumps covering the data changes in three months and divide the data into 300 versions. Each version has an average size of 855MB. We generate the read and write workload using keys uniformly selected from the dataset to test the throughput.

*5.1.3 Ethereum Transactions.* We use real-world Ethereum transaction data[3] from Block 8900000 to 9200000, where the key is the 64-bytes block hash and the value is the RLP (Recursive Length Prefix) encoded raw transaction data. The length of the raw transaction ranges from 100 bytes to 57738 bytes with an average of 532 bytes. RLP is the main encoding method used to serialize objects in Ethereum, which is also used to encode raw transactions. In Ethereum, each block naturally makes a new version.

## 5.2 Implementation

In this section, we briefly describe the implementation of selected indexes and the baseline. We port the Ethereum's implementation [3] of MPT to our experiment environment, which adopts the path compaction optimization. The implementation of MBT is based on the source code provided in Hyperledger Fabric 0.6 [7]. We further make it immutable and add index lookup logic, which is missing in the original implementation. For POS-Tree, we use the implementation in Forkbase [44]. Moreover, we further apply batching techniques, taking advantage of the bottom-up build order, to reduce the number of tree traversal and hash calculations significantly. Lastly, to compare *SIRI* and non-*SIRI* structures, we implement an immutable B$^+$-tree with tamper evidence support, called Multi-Version Merkle B$^+$-tree (MVMB$^+$-Tree), as the baseline. We replace the pointers stored in index nodes with the hash of their immediate children and maintain an additional table from the hash to the actual address. For all the structures, we adopt node-level copy-on-write to achieve the data immutability.

## 5.3 Throughput and Latency

We evaluate the three candidates and the baseline from a traditional view in this part, where throughput and latency are the major measurements.

*5.3.1 Throughput.*
First, we evaluate the throughput using the YCSB dataset. We run the read, the write and the mixed workloads under diverse data size and skewness. The results are illustrated in Figure 6. It can be observed that the throughput of all indexes decreases as the number of data grows and complies with the operation bound formulated in Section 4.1. Figure 6(a) shows the throughput for the read workload with uniform data. The throughput of POS-Tree is 0.95x - 1.06x of the baseline while MPT is only 0.74x - 0.96x of the baseline. The throughput of MBT drops quickly from 3.2x to 0.45 of the baseline due to the dominating leaf loading and scanning process. As it is shown in Figure 13, the time to traverse the tree and load the nodes keeps constant, while time to scan leaf node keeps increasing. For the write workload shown in Figure 6(c), we can observe a similar trend. However, POS-Tree performs 1.04x-1.9x better than the baseline taking advantage of the batching techniques and the bottom-up building process.

By comparing Figure 6 horizontally, we can observe that the throughput of all data structures decreases drastically as the ratio of write operations increases. This is due to the cost of node creation, memory copy and cryptographic function
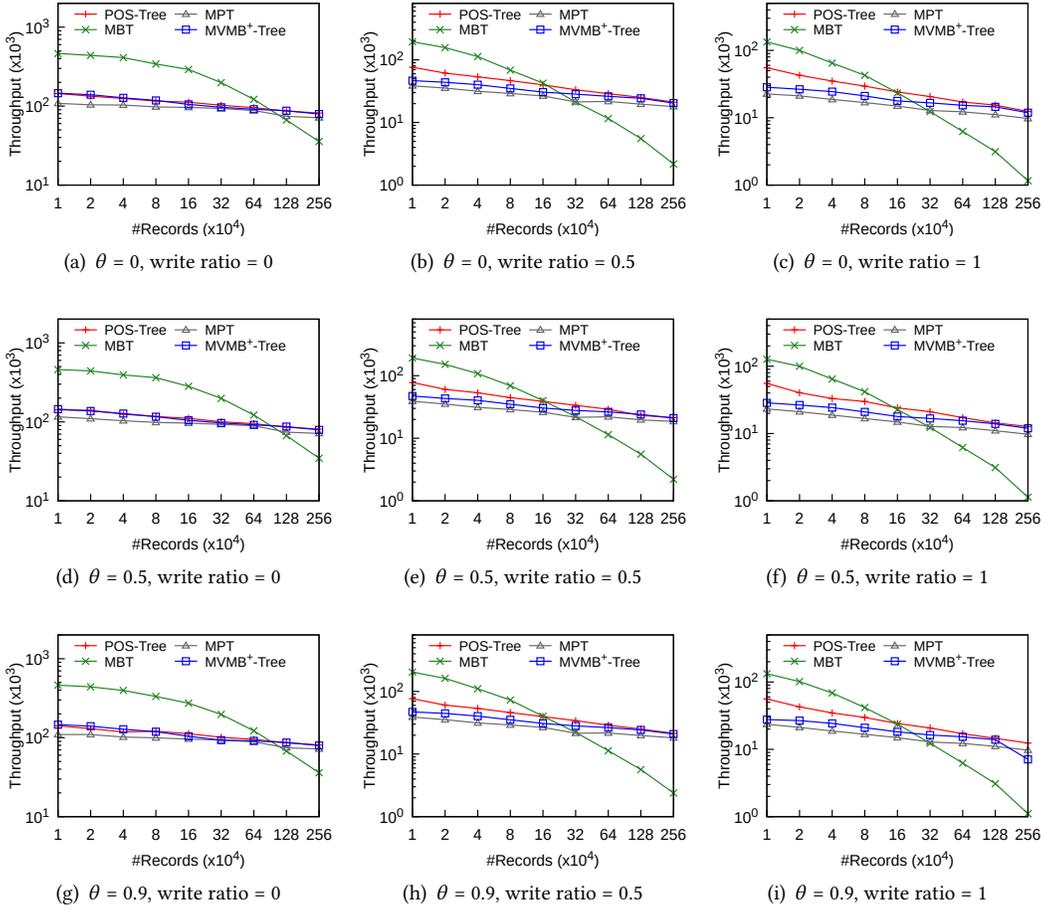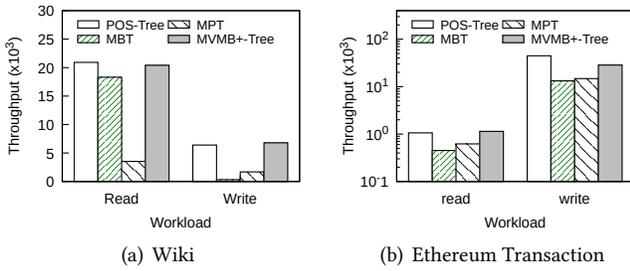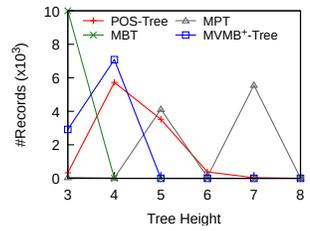
---

[2]https://dumps.wikimedia.org/enwiki/

[3]https://cloud.google.com/blog/products/data-analytics/ethereum-bigquery-public-dataset-smart-contract-analytics

(a) $\theta = 0$, write ratio = 0

(b) $\theta = 0$, write ratio = 0.5

(c) $\theta = 0$, write ratio = 1

(d) $\theta = 0.5$, write ratio = 0

(e) $\theta = 0.5$, write ratio = 0.5

(f) $\theta = 0.5$, write ratio = 1

(g) $\theta = 0.9$, write ratio = 0

(h) $\theta = 0.9$, write ratio = 0.5

(i) $\theta = 0.9$, write ratio = 1

Figure 6: Throughput on YCSB



(a) Wiki

(b) Ethereum Transaction

Figure 7: Throughput on real world datasets

Figure 8: Diff performance

Figure 9: Tree height

computation. The absolute throughput drops over 6.6x in the comparison of the largest dataset for POS-Tree and baseline while it drops 30x for MBT and 7.3x for MPT. By comparing Figure 6 vertically, we can observe that there is no change in throughput for all index structures when $\theta$ changes from 0 to 0.9. Therefore, we can conclude that they are all resilient to data skewness.

It is also worth noting that, compared with MBT, the other three structures perform much more steady for both read and write workloads. Meanwhile, POS-Tree outperforms MPT in all cases and has comparable performance compared to our baseline index.

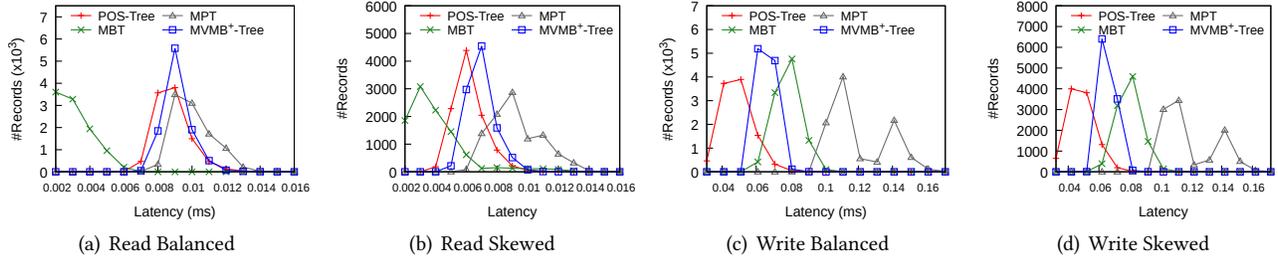Next, we run the experiment on the Wiki dataset. The system first load the entire dataset batched in 300 versions,

(a) Read Balanced     (b) Read Skewed     (c) Write Balanced     (d) Write Skewed

**Figure 10: Latency on YCSB**



(a) Read     (b) Write         (a) Read     (b) Write
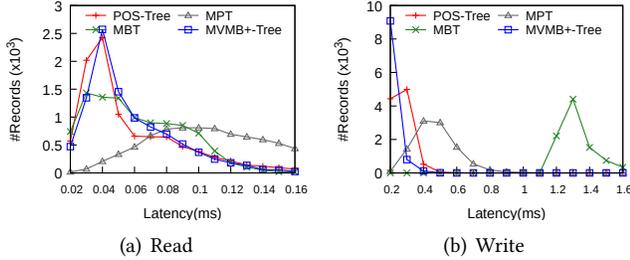
**Figure 11: Latency on Wiki data**      **Figure 12: Latency on Ethereum transaction data**
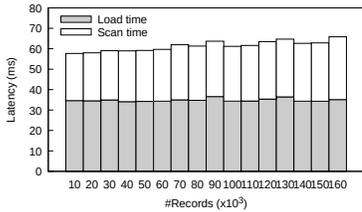


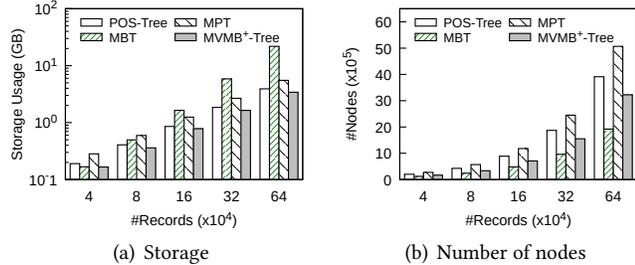**Figure 13: MBT breakdown latency**

(a) Storage     (b) Number of nodes

**Figure 14: Performance on single group data access**

and then execute the read and write workloads that is uniformly selected. Figure 7(a) demonstrates the results that are aligned with those in the YCSB experiment.

Lastly, for the experiments on Ethereum data, we simulate the way blockchain stores the transactions. For each block, we build an index on transaction hash for all transactions within that block and store the root hash of the tree in a global linked list. Versions are naturally created at a block granularity. For write operations, the system appends the new block of transactions to the global linked list while for lookup operations, it scans the linked list for the block containing the transaction, and traverses the index to obtain the value. Figure 7(b) shows the result of this experiment. It can be observed that POS-Tree outperforms other indexes in write workloads. This is because we are building indexes for each block instead of a global index. Further, instead of insert/update operations, we perform batch loading from

scratch. In this case, POS-Tree's bottom-up building process is superior to the MPT's and MVMB$^+$-Tree's top-down building process, as it only traverses the tree and creates each node once. Another difference is that the throughput of read workload is lower than that of the write workload mainly due to the additional block scanning time.

### 5.3.2 Latency and Path Length.

In this experiment, we measure the latency of each read and write operation and calculate the distribution with balanced and skewed data. For the YCSB dataset, read-only and write-only workloads are fed into the indexes with balanced ($\theta = 0$) and highly skewed ($\theta = 0.9$) distributions. The dataset used in this test contains 160,000 keys. We run 10,000 operations and pictured the latency distribution in Figure 10. The x-axis is the range of the latency and the y-axis is the number of records fell in that latency range. It can be seen from the figure that the rankings among the indexes coincide

with the previous conclusion – POS-Tree performs the best for both read and write workloads while MPT performs the worst. Meanwhile, MPT has several peak points, representing operations accessing data stored in different levels of the tree. MBT experiences the most dramatic changes between read and write workloads. It outperforms all the candidates in the read workloads but is worse than POS-Tree in write workload.

To take a closer observation of how the workloads affect the candidates, we further gather the traversed tree height of each operation for the write-only workload with the uniform distribution. The results are shown in Figure 9(a), where the x-axis represents the height of the lookup path and the y-axis indicates the number of operations. Most operations have to visit 4-level nodes to reach the bottom-most level of POS-Tree whilst 5- or 7-level nodes are frequently traversed for MPT. The efficiency in MBT is also verified in the figure since all requests only need 3 levels to reach the bottom of the structure in both balanced and skewed scenarios.

We can obtain similar results and conclusions on the Wiki dataset as shown in Figure 11. However, the experiment on Ethereum transaction data exhibits different trends as depicted in Figure 12. As can be observed, all the structures have similar read latency, caused by the dominant block scanning process.

We also run a diff workload to evaluate the performance of "diff" operations. In the experiment, each structure loads two versions of data in random order. A diff operation is performed between the two versions and the execution time is taken, as depicted in Figure 8(a). All the candidates outperform the baseline due to the structurally invariant property. Among which, MBT performs the best (4x of baseline) since the position of the nodes containing a specific data is static among all versions. The logic of diff operation is the simplest, i.e., comparing the hash of the nodes at the corresponding position. MPT performs 2x better than the baseline and 1.7x better than POS-Tree due to the simplicity that keys with the same length always lie in the same level of the tree.

## 5.4 Storage

In this section, we evaluate the space consumption of the index structures under different use cases.

### 5.4.1 Single Group Data Access.

We first start with a simple case, where a dataset is accessed by multiple users. There is no sharing of data or cross-department collaborative editing in this setting. Therefore, the deduplication benefit is limited using *SIRI*. In reality, such case often happens in-house within a single group of users from the same department. Figure 14(a) shows the storage under different data sizes for the YCSB dataset. There are two main factors affecting the space efficiency, i.e., the

size of the node and the height of the tree. On the one hand, larger tree height results in more node creations for write operations, which also increases the space consumption. As an example, MPT performs badly since it has the largest tree height in our experiment setting. It consumes the storage up to 1.6x higher than the baseline and up to 1.4x larger than POS-Tree. On the other hand, a large node size means that even minor changes to the node could trigger the creation of a new substantial node, which hence leads to larger space consumption. As can be seen, MBT performs the worst due to the largest node size it has in the implementation. It consumes up to 6.4x the space of that used by the baseline. POS-Tree, compared to the baseline, also has a larger node size variance due to content-defined chunking, leading to a greater number of large nodes.

To better analyze how the memory space is used by different pages, we further accumulate the number of nodes for all chosen indexes. The results are demonstrated in Figure 14(b) with variant dataset sizes. Typically, they follow similar trends as Figure 14(a), except that MBT generates the least number of nodes as the total number of nodes is fixed for the structure. The reason is rooted from the nature of MBT, which has a fixed total number of nodes and increasing leaf node size as more records are inserted. Therefore, the number of nodes created keeps constant when updating or inserting, no matter how large the total number of records is. On the contrary, other structures have a fixed node size and an increasing number of nodes, causing the number of nodes created, as well as the height of the tree, increases during updating or inserting.

The results for the Wiki dataset and Ethereum transaction dataset are shown in Figure 15 and Figure 16. Similar to the results of the YCSB experiment, MBT and MPT consumed more space than POS-Tree and MVMB$^+$-Tree. A difference is that MPT storage consumption increases very fast as the number of versions are loaded. This is because the key length of the Wiki dataset is much larger than that of YCSB, and the encoding method used by Ethereum further doubles the key length. This makes MPT a very sparse tree. For every insert/update operation, more nodes need to be re-hashed and created. Hence, the space efficiency is worse than it shows in the YCSB experiment.

Another difference is that MBT generates more nodes compared with other experiments. This is again because of the experiment setting that a new instance of index will be created per block. Since each block only contains a few hundreds of transactions, MBT is less efficient compared with other structures.

### 5.4.2 Diverse Group Collaboration.

We now examine the storage consumption in applications where different groups of users are collaborating to work on
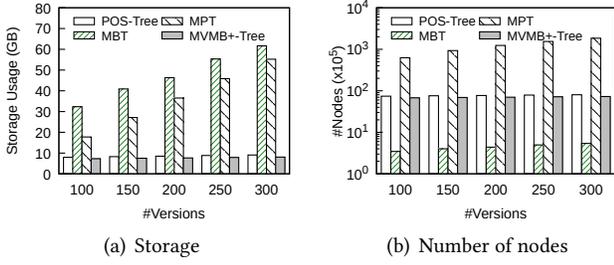
(a) Storage      (b) Number of nodes
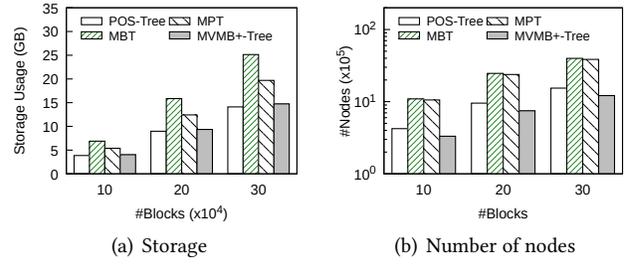
**Figure 15: Storage on Wiki data**



(a) Storage      (b) Number of nodes

**Figure 16: Storage on Ethereum transaction data**



(a) Storage    (b) Number of nodes    (c) Deduplication ratio    (d) Node sharing ratio

**Figure 17: Performance on diverse group collaboration with varying overlap ratio**



(a) Storage    (b) Number of nodes    (c) Deduplication ratio    (d) Node sharing ratio

**Figure 18: Performance on diverse group collaboration with varying batch size**

the same dataset. This often occurs in the data cleansing process and data analysis procedure, where diverse parties work on different parts of the same dataset. One significant phenomenon in this use case is that duplicates can be frequently found. Therefore, the deduplication capability introduced by *SIRI* is critical to improving the space efficiency. To evaluate the deduplication capability, we define another metric called node sharing ratio from a different aspect. The metric can be formulated as follows:

$$\eta(S) = 1 - \frac{|P_1 \cup P_2 \cup ... \cup P_k|}{|P_1| + |P_2| + ... + |P_k|},$$

where $P_i$ is the set of nodes of an instance i. While the deduplication ratio evaluates the size of the storage saved, the node sharing ratio indicates how many duplicate nodes have been eliminated.

The YCSB dataset is used in this experiment. We simulate 10 groups of users, each of which initializes the same dataset of 40,000 records. We generate workloads of 160,000 records with overlap ratios ranging from 10% to 100% and feed them to the candidates. Here, 10% overlap ratio means 10% of the records have the same key and value. The execution is processed with default batch size, i.e., 4,000 records.

The results of the deduplication ratio and the node sharing ratio are shown in Figure 17(c) and Figure 17(d), respectively. Both metrics of all the structures become higher when the workload overlap ratio increases since more duplicate nodes can be found due to increasing similarities among the datasets. Benefiting from smaller node size and smaller portion of updating nodes, MPT achieves the highest deduplication ratio (up to 0.96) and node sharing ratio (up to 0.7).

13

POS-Tree achieves a slightly better deduplication ratio than the baseline though they both have similar size of nodes and the height of the tree. The actual ratios are 0.88 and 0.86, respectively. However, it achieves a much better node sharing ratio compared to the baseline (0.48 vs. 0.27) because of its content-addressable strategy when chunking the data. By contrast, MBT's fixed number of pages and growing leaf nodes limit the number of duplicates, and therefore it does not perform as good as the other two *SIRI* representatives.

To be more precise, we further collect the storage usage and the number of pages created by the testing candidate and illustrate the results in Figure 17(a) and Figure 17(b). The trends in the figures match the corresponding deduplication ratio and node sharing ratio perfectly. With the increasing overlap ratio, storage reduction of POS-Tree and MPT is more obvious than the baseline, among which POS-Tree is the most space-efficient. MPT is most sensitive to overlap ratio changes due to the high node sharing ratio introduced by its structural design. Although it consumes more space for non-overlapping datasets, MPT outperforms the baseline for datasets with the overlap ratio above 90%.

We also evaluate the effect of query batch size on storage space. Same as previous experiments' setting, we simulate 10 parties. Each of them initializes the same dataset contains 40,000 records and executes workloads containing 160,000 keys with default overlap ratio 50%. Figure 18(c) depicts how the deduplication ratio decreases along with the query batch size increases. The reason is that larger batch sizes cause a larger portion of the index structure to be updated, resulting in fewer nodes to be reused between versions. Figure 18(a) and Figure 18(b) show the storage usage and Number of nodes created with different batch sizes. Similar relationships across the structures can be observed as in Figure 17(a) and Figure 17(b). Except for both of the metrics decrease when using a larger batch size due to less versions stored in the index.

### 5.4.3 Structure Parameters.
The parameters of the indexes can affect the deduplication ratio as aforementioned in section 4. The impact of those key parameters is verified in this experiment, namely node size for POS-Tree, number of buckets for MBT and mean key length for MPT. For POS-Tree, the boundary pattern is varied to change the node size from 512 to 4,096 bytes probabilistically. For MBT, the number of fixed buckets is set from 4000 to 10,000. For MPT, the dataset is generated with different minimum key lengths, which can lead to diverse mean key length from 10.2 to 13.7. (The maximum key length in the dataset is fixed.) The results are shown in Table 3, which coincide with the conclusions in Section 4. The deduplication ratio of POS-Tree increases as the average node size increases. This is expected as the number of same

large nodes is less than that of small nodes, leading to fewer occurrences of duplicate pages. Similarly, the deduplication ratio of MBT increases as the number of buckets increases because a larger number of buckets results in smaller leaf nodes. The deduplication ratio of MPT increases as the mean key length increases. This is because longer keys usually have more conflicting bits and result in a wider tree. Therefore, the portion of the reusable nodes increases.

## 5.5 Breakdown Analysis

In this section, we evaluate how each *SIRI* property affects the storage and deduplication performance. We select POS-Tree as the testing object and disable the properties one by one. For each property, we first explain how each property is disabled and then provide the experimental results following closely. We note that the *Universally Reusable* property is common for all immutable tree indexes using copy-on-write approach. Thus, it is ignored in this experiment.

### 5.5.1 Disabling Structurally Invariant Property.
The pattern-aware partitioning is the key to guarantee the *Structurally Invariant* property. Therefore, we disable the property by forcibly splitting the entries at half of the maximum size when no pattern is found within the maximum size. Consequently, the resulting structure depends on the data insertion order. We increase the probability of not finding the pattern by increasing the bits of pattern and lowering the maximum value.

The result is presented in Figure 19(a). We can observe an up to 15% decrease in the deduplication ratio when *Structurally Invariant* property is disabled. For instance, the deduplication ratio drops from 0.67 to 0.52 when the workload overlap ratio equals to 100%. It is expected as the index performs the operations in different orders, resulting in different nodes and smaller number of share-able pages. Though the records stored are the same, POS-Tree cannot reuse the nodes with the *Structurally Invariant* property disabled. Similarly, Figure 19(b) shows that the node sharing ratio decreases by up to 17%, i.e. from 0.53 to 0.36, by disabling *Structurally Invariant* property.
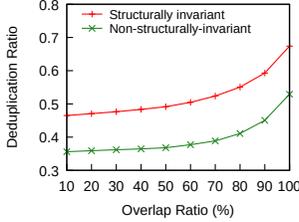
### 5.5.2 Disabling Recursively Identical Property.
Originally, only the set of nodes lying in the path from the root to the leaf node is copied and modified when an update operation is performed, while the rest of the nodes are shared between the two versions in POS-Tree. We disable *Recursively Identical* property by forcibly copying all nodes in the tree. The number of different pages between the two instances is much larger than the number of intersections, which is zero.
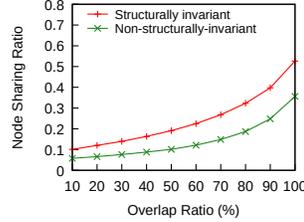
Figure 20(a) shows that the deduplication ratio for POS-Tree with *Recursively Identical* disabled is 0 since the structure does not allow the sharing of nodes among different

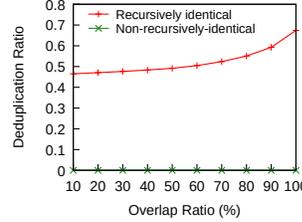**Table 3: Effect of structure parameters on the deduplication ratio**

| Node Size | $\eta$(POS-Tree) | | # Buckets | $\eta$(MBT) | | $\overline{keylen}$ | $\eta$(MPT) |
|---|---|---|---|---|---|---|---|
| 512 | 0.722 | | 4000 | 0.3301 | | 10.2 | 0.9685 |
| 1024 | 0.6485 | | 6000 | 0.4599 | | 12 | 0.9693 |
| 2048 | 0.5391 | | 8000 | 0.5433 | | 13.3 | 0.9806 |
| 4096 | 0.4108 | | 10000 | 0.6003 | | 13.7 | 0.9823 |



(a) Deduplication ratio  (b) Node sharing ratio

**Figure 19: Effect of Structurally Invariant property**



(a) Deduplication ratio  (b) Node sharing ratio

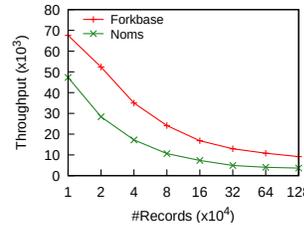**Figure 20: Effect of Recursively Identical property**
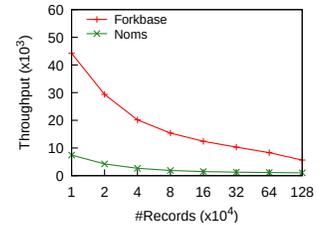


(a) Read  (b) Write

**Figure 21: Performance integrated with Forkbase**



(a) Read  (b) Write

**Figure 22: Comparison between Forkbase and Noms**

versions. Obviously, the node sharing ratio of non *Recursively Identical* POS-Tree shown in Figure 20(b) is also 0. Compared to the figures in previous sections, we can infer how this property accelerates the deduplication rate and ultimately influences the final storage performance.

Overall, we can conclude that *Recursively Identical* property is the fundamental property to enable indexes with deduplication and node sharing across different users and datasets. On top of this, *Structurally Invariant* property further enhances the level of deduplication and node sharing by making structures history-independent.

## 5.6 System Experiment

### 5.6.1 Integration with Forkbase.

To further evaluate the performance of *SIRI*, we integrate the indexes into Forkbase [44], a storage engine for blockchain and forkable applications. In this experiment, we configured a single Forkbase servlet and a single client

to benchmark the system-level throughput. The evaluation results are demonstrated in Figure 21.

For read operations, the main difference between index-level performance and system integrated performance is the remote access due to client-server architecture. The overhead of remote access becomes the dominant factor of performance. To mitigate such overhead, Forkbase caches the nodes at clients after retrieved from servers. Hence, the following read operations on the same nodes can benefit from performing only local access.

Figure 21(a) shows the throughput of read workload. Similar to index-level experiments, the throughput decreases when the total number of records grows. POS-Tree achieves comparable performance to our baseline MVMB+-Tree, and it outperforms the other 2 indexes when the total number of records is large (greater than 2,560,000). MPT performs the worst among all indexes due to larger tree height, which comply with the operation bound in Section 4.1.1. Different from

the index-level experiment, MBT performs worse than POS-Tree and MVMB+-Tree when the number of records is extremely small (10,000 records). This is because the hit ratio of cached nodes for MBT is lower than other indexes. Since all index nodes of MBT have a fixed number of entries, the number of repeated reads is less compared with POS-Tree and MVMB+-Tree, where large nodes contribute more repeated reads. When the number of records grows larger, the affected portion of POS-Tree and MVMB+-Tree decreases. Consequently, the number of repeated reads decreases. While the structure of MVMB+-Tree keeps unchanged, leading to a constant number of repeated reads. Therefore, MBT performs better when the number of records is greater than 20,000. When the number of records is greater than 2,560,000, the bottleneck becomes the loading time and the scanning time of leaf nodes, and the throughput drops below that of other indexes.

The write operations will be performed on the server side completely. Hence they will not be affected by the hit ratio of cached nodes described above. Figure 21(b) shows the throughput of write workload. We can observe similar results as that of index-level experiments.

### 5.6.2 Comparison between Forkbase and Noms.

Next, we perform a comparative study between Forkbase and Noms [10]. Both systems facilitate data versioning management using similar indexing concept. As in POS-Tree, the bottom most layer of Noms' Prolly Tree uses the sliding-window approach to partition the leaf nodes based on the boundary pattern. To match the boundary pattern in the internal layers, POS-Tree directly uses hash values of the child nodes, while Prolly Tree uses the hash values repeatedly computed from the sliding-window. Such computational overhead causes inefficiency of its write operations.

In the experiment, we directly use the code of Noms from its Github repository, which is implemented in GO. We use Noms' remote setup on top of their own HTTP protocol to compare with Forkbase's single server single client setup as described previously. To make a fair comparison, we configure the node size of POS-Tree to 4K with window size of 67 bytes, which is the default setting of Noms. The experiment is conducted as follows. First, we initialize the systems with 10K to 128K records. Then we execute read and write workload of 10K records respectively to measure the throughput. The results are shown in Figure 22. We can observe that Forkbase performs 1.4x-2.7x better in read operations and 5.6x-8.4x better in write operations than Noms.

## 6 CONCLUSION

Tamper evidence and deduplication are two properties increasingly demanded in emerging applications on immutable data, such as digital banking, blockchain and collaborative

analytics. Recent works [7, 44, 45] have proposed three index structures equipped with these two properties. However, there have been no systematic comparisons among them. To address the problem, we conduct a comprehensive analysis of all three indexes in terms of both theoretical bounds and empirical performance. Our analysis provides insights regarding the pros and cons of each index, based on which we conclude that POS-Tree [44] is a favorable choice for indexing immutable data.

## REFERENCES

[1] Amazon quantum ledger database. https://aws.amazon.com/qldb/.
[2] Azure blockchain service. https://azure.microsoft.com/en-us/services/blockchain-service/.
[3] Ethereum. https://www.ethereum.org.
[4] Git. https://git-scm.com/.
[5] Git(and github) for data. https://blog.okfn.org/2013/07/02/git-and-github-for-data/.
[6] Googledocs. https://www.docs.google.com.
[7] Hyperledger. https://www.hyperledger.org.
[8] LevelDB. https://github.com/google/leveldb.
[9] Mercurial. https://www.mercurial-scm.org/.
[10] Noms. https://github.com/attic-labs/noms.
[11] RocksDB. http://rocksdb.org.
[12] Subversion. https://subversion.apache.org/.
[13] Webank. https://www.webank.com/en/.
[14] I. Ahn and R. Snodgrass. Performance evaluation of a temporal database management system. In *SIGMOD Record*, volume 15, pages 96–107, 1986.
[15] A. Arasu, K. Eguro, R. Kaushik, D. Kossmann, P. Meng, V. Pandey, and R. Ramamurthy. Concerto: A high concurrency key-value store with integrity. In *SIGMOD*, 2017.
[16] H. Berenson, P. Bernstein, J. Gray, J. Melton, E. O'Neil, and P. O'Neil. A critique of ansi sql isolation levels. In *SIGMOD*, pages 1–10, 1995.
[17] P. A. Bernstein and N. Goodman. Multiversion concurrency control-theory and algorithms. *ACM Trans. Database Syst.*, 8(4):465–483, Dec. 1983.
[18] A. Bhardwaj, S. Bhattacherjee, A. Chavan, A. Deshpande, A. J. Elmore, S. Madden, and A. Parameswaran. Datahub: Collaborative data science & dataset version mangement at scale. In *CIDR*, 2015.
[19] T. Chiueh and D. Pilania. Design, implementation, and evaluation of a repairable database management system. In *ICDE*, 2005.
[20] T. T. A. Dinh, R. Liu, M. Zhang, G. Chen, B. C. Ooi, and J. Wang. Untangling blockchain: A data processing view of blockchain systems. *TKDE*, 30(7):1366–1385, 2018.
[21] T. T. A. Dinh, J. Wang, G. Chen, R. Liu, B. C. Ooi, and K.-L. Tan. Blockbench: A framework for analyzing private blockchains. In *SIGMOD*, pages 1085–1100, 2017.
[22] K. Eshghi and H. K. Tang. A framework for analyzing and improving content-based chunking algorithms. *Hewlett-Packard Labs Technical Report TR*, 30(2005), 2005.
[23] S. Huang, L. Xu, J. Liu, A. J. Elmore, and A. Parameswaran. OrpheusDB: Bolt-on versioning for relational databases. *PVLDB*, 10(10):1130–1141, 2017.

[24] G. Kollios and V. J. Tsotras. Hashing methods for temporal data. *IEEE Trans. on Knowl. and Data Eng.*, 14(4):902–919, 2002.

[25] J. Kubiatowicz, D. Bindel, Y. Chen, S. Czerwinski, P. Eaton, D. Geels, R. Gummadi, S. Rhea, H. Weatherspoon, W. Weimer, C. Wells, and B. Zhao. Oceanstore: An architecture for global-scale persistent storage. In *ASPLOS*, 2000.

[26] S. Lanka and E. Mays. Fully persistent B+-trees. In *SIGMOD*, pages 426–435, 1991.

[27] M. Maddox, D. Goehring, A. J. Elmore, S. Madden, A. G. Parameswaran, and A. Deshpande. Decibel: The relational dataset branching system. *PVLDB*, 9(9):624–635, 2016.

[28] R. C. Merkle. A digital signature based on a conventional encryption function. In *A Conference on the Theory and Applications of Cryptographic Techniques on Advances in Cryptology*, pages 369–378, 1988.

[29] A. Muthitacharoen, B. Chen, and D. Mazieres. A low-bandwidth network file system. In *SIGOPS Operating Systems Review*, volume 35, pages 174–187, 2001.

[30] S. Nakamoto. Bitcoin: A peer-to-peer electronic cash system. https://bitcoin.org/bitcoin.pdf, 2009.

[31] C. Okasaki. Purely functional data structures. *Cambridge University Press*, 1999.

[32] J. Paulo and J. Pereira. A survey and classification of storage deduplication systems. *ACM Computing Surveys (CSUR)*, 47(1):11, 2014.

[33] D. R. K. Ports and K. Grittner. Serializable snapshot isolation in postgresql. *PVLDB*, 5(12):1850–1861, 2012.

[34] J. R. Driscoll, N. Sarnak, D. Sleator, and R. Tarjan. Making data structures persistent. volume 38, pages 109–121, 01 1989.

[35] O. Rodeh. B-trees, shadowing, and clones. *ACM Transactions on Storage (TOS)*, 3(4), 2008.

[36] P. Ruan, G. Chen, T. T. A. Dinh, Q. Lin, B. C. Ooi, and M. Zhang. Fine-grained, secure and efficient data provenance on blockchain systems. *PVLDB*, 12(9):975–988, 2019.

[37] B. Salzberg and V. J. Tsotras. Comparison of access methods for time-evolving data. *ACM Computing Surveys (CSUR)*, 31(2):158–221, 1999.

[38] D. J. Santry, M. J. Feeley, N. C. Hutchinson, and A. C. Veitch. Elephant: the file system that never forgets. In *HotOS*, pages 2–7, 1999.

[39] A. Seering, P. Cudre-Mauroux, S. Madden, and M. Stonebraker. Efficient versioning for scientific array databases. In *ICDE*, pages 1013–1024, April 2012.

[40] C. A. N. Soules, G. R. Goodson, J. D. Strunk, and G. R. Ganger. Metadata efficiency in versioning file systems. In *FAST*, 2003.

[41] M. Stonebraker and L. A. Rowe. The design of the POSTGRES. In *SIGMOD*, pages 340–355, 1986.

[42] J. D. Strunk, G. R. Goodson, M. L. Scheinholtz, C. A. N. Soules, and G. R. Ganger. Self-securing storage: Protecting data in compromised system. In *OSDI*, 2000.

[43] A. U. Tansel, J. Clifford, S. Gadia, S. Jajodia, A. Segev, and R. Snodgrass. Temporal databases: Theory, design, and implementation. *Benjamin-Cummings Publishing Co., Inc.*, 1993.

[44] S. Wang, T. T. A. Dinh, Q. Lin, Z. Xie, M. Zhang, Q. Cai, G. Chen, B. C. Ooi, and P. Ruan. Forkbase: An efficient storage engine for blockchain and forkable applications. *PVLDB*, 11(10):1137–1150, 2018.

[45] D. Wood. Ethereum: A secure decentralised generalised transaction ledger. 2014.

[46] W. Xia, H. Jiang, D. Feng, F. Douglis, P. Shilane, Y. Hua, M. Fu, Y. Zhang, and Y. Zhou. A comprehensive study of the past, present, and future of data deduplication. *Proceedings of the IEEE*, 104(9):1681–1710, 2016.

[47] N. Zhu and T.-C. Chiueh. Design, implementation, and evaluation of repairable file service. In *DSN*, 2003.