



Hybrid Blockchain Database Systems: Design and Performance

Zerui Ge
School of Computing
National University of Singapore
gezerui@comp.nus.edu.sg

Dumitrel Loghin
School of Computing
National University of Singapore
dumitrel@comp.nus.edu.sg

Beng Chin Ooi
School of Computing
National University of Singapore
ooibc@comp.nus.edu.sg

Pingcheng Ruan
School of Computing
National University of Singapore
ruanpc@comp.nus.edu.sg

Tianwen Wang
School of Computing
National University of Singapore
tianwenw@nus.edu.sg

ABSTRACT

With the emergence of hybrid blockchain database systems, we aim to provide an in-depth analysis of the performance and trade-offs among a few representative systems. To achieve this goal, we implement Veritas and BlockchainDB from scratch. For Veritas, we provide two flavors to target the crash fault-tolerant (CFT) and Byzantine fault-tolerant (BFT) application scenarios. Specifically, we implement Veritas with Apache Kafka to target CFT application scenarios, and Veritas with Tendermint to target BFT application scenarios. We compare these three systems with the existing open-source implementation of BigchainDB. BigchainDB uses Tendermint for consensus and provides two flavors: a default implementation with blockchain pipelining and an optimized version that includes blockchain pipelining and parallel transaction validation. Our experimental analysis confirms that CFT designs, which are typically used by distributed databases, exhibit much higher performance than BFT designs, which are specific to blockchains. On the other hand, our extensive analysis highlights the variety of design choices faced by the developers and sheds some light on the trade-offs that need to be done when designing a hybrid blockchain database system.

PVLDB Reference Format:

Zerui Ge, Dumitrel Loghin, Beng Chin Ooi, Pingcheng Ruan, and Tianwen Wang. Hybrid Blockchain Database Systems: Design and Performance. PVLDB, 15(5): 1092 - 1104, 2022.
doi:10.14778/3510397.3510406

PVLDB Artifact Availability:

The source code, data, and/or other artifacts have been made available at <https://github.com/nusdbsystem/Hybrid-Blockchain-Database-Systems>.

1 INTRODUCTION

In the last few years, a handful of systems that integrate both distributed databases and blockchain properties have emerged in the academic database community [16, 19, 20, 31, 33, 38]. These systems,

termed *hybrid blockchain database systems*, are either adding database features to existing blockchains to improve performance and usability or implement blockchain features in distributed databases to improve their security [36]. However, there is little comparison among these systems. In this paper, we aim to fill this gap by providing an in-depth analysis of a few representative hybrid blockchain database systems.

To achieve this goal, we first have to overcome a challenge: only one such system, namely BigchainDB [20] is open-source. Moreover, by the time of writing this paper, we did not manage to get the source code of any other hybrid blockchain database system. Hence, we undertake the tedious task of implementing Veritas [19] and BlockchainDB [16]. By doing so, we gain the flexibility of changing some parts of the systems to better compare them to other systems. For example, we can change the underlying database from a relational SQL type to a NoSQL type. Or we can change the consensus mechanism from crash fault-tolerant (CFT) to Byzantine fault-tolerant (BFT).

The original design of Veritas uses Apache Kafka [17], which is a CFT service, as the broadcasting service among server nodes. That is, when a server node needs to update its local shared database as a result of a transaction's execution, it sends this update to all the other server nodes via the Kafka service. Moreover, the other server nodes send their agreement or disagreement regarding that update via the Kafka service. Apache Kafka uses a primary backup mechanism to achieve CFT. Hence, it has high efficiency at the cost of decreased liveness. On the other hand, most blockchain systems adopt a BFT consensus mechanism. Hence, we also implement a version of Veritas where Tendermint [3], which is a BFT consensus, is used as middleware, similar to BigchainDB [20].

BlockchainDB [16] implements a shared (and sharded) database on top of a classic blockchain. The database interface provides a simple key-value store API to the user, similar to Veritas and BigchainDB. The blockchain layer is flexible, being able to interact with different blockchains, such as Ethereum [4] and Hyperledger Fabric [18]. In this paper, as well as in [16], BlockchainDB uses Ethereum as its underlying blockchain.

BigchainDB introduces two optimizations. The first optimization, called blockchain pipelining, allows nodes to vote for a new block while the current block of the ledger is still undecided. Each node will just reference its last decided block when voting for a new block. By doing so, the system avoids waiting for blocks to be fully committed before proposing new blocks. Second, BigchainDB

This work is licensed under the Creative Commons BY-NC-ND 4.0 International License. Visit <https://creativecommons.org/licenses/by-nc-nd/4.0/> to view a copy of this license. For any use beyond those covered by this license, obtain permission by emailing info@vldb.org. Copyright is held by the owner/author(s). Publication rights licensed to the VLDB Endowment.
Proceedings of the VLDB Endowment, Vol. 15, No. 5 ISSN 2150-8097.
doi:10.14778/3510397.3510406

includes a flavor called BigchainDB with parallel validation, or BigchainDB (PV), where transactions are validated in parallel on multiple CPU cores. These parallel validations should increase the throughput, however, we do not observe any improvement, as we shall see in our experimental analysis.

In summary, we make the following contributions:

- We survey and qualitatively compare the existing hybrid blockchain database systems.
- We provide flexible implementations of BlockchainDB and Veritas. In particular, we implement Veritas with Apache Kafka and Tendermint as mechanisms to broadcast the transactions. While Kafka provides crash fault tolerance, Tendermint is designed to work in a Byzantine environment, closer to what the blockchains are supposed to do.
- We analyze the performance of five systems, namely, Veritas (Kafka), Veritas (Tendermint), BlockchainDB, BigchainDB, and BigchainDB with parallel validation. Our analysis exposes the trade-offs to be considered by the developers to achieve the best performance in their application scenario.
- Among others, we show that Veritas (Kafka) exhibits a higher performance compared to all the other systems. For example, Veritas (Kafka) exhibits close to 30,000 TPS, while Veritas (TM), BigchainDB, and BigchainDB (PV) exhibit close to 1,700, 180, and 180 TPS, respectively, on networks of four peers. On the other hand, BlockchainDB exhibits less than 100 TPS. While there is room for optimization in all the systems, the significant gap in performance between CFT and BFT consensus mechanisms will be hard to reduce.

The rest of this paper is organized as follows. In Section 2 we survey the existing hybrid blockchain database systems. In Section 3, we first describe our implementations of Veritas and BlockchainDB, followed by the existing open-source implementation of BigchainDB. In Section 4 we conduct our performance analysis. In Section 5 we discuss the trade-offs and challenges that users and developers encounter when using a hybrid blockchain database system. Finally, we conclude the paper in Section 6.

2 BACKGROUND AND RELATED WORK

In the last few years, there have been a few works published in database conferences that integrate blockchains and databases [16, 19, 20, 31, 33, 38]. Such systems allow companies to do business with peace of mind because every database operation is tracked by a distributed ledger. However, different business scenarios have different requirements, and as a result, many hybrid blockchain database systems have been built for different application scenarios. In this section, we describe the existing *hybrid blockchain database systems* and we briefly mention some similar systems that can be classified as *ledger databases*.

2.1 Hybrid Blockchain Database Systems

Veritas [19] is a shared database design that integrates an underlying blockchain (ledger) to keep auditable and verifiable proofs. The interaction between the database and the blockchain is done through verifiers. A verifier takes the transaction logs from the database, makes a correctness-checking decision, and sends it to the other verifiers. The final decision agreed by all the verifiers is

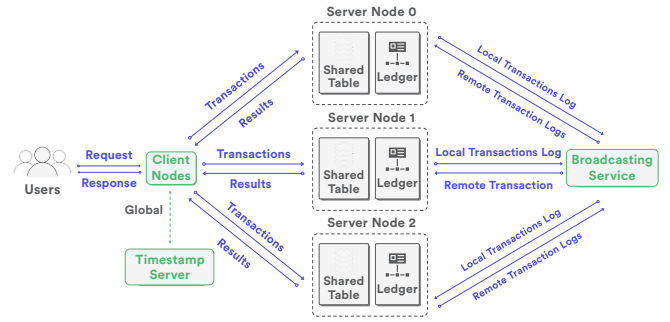


Figure 1: Veritas (Kafka) with Apache Kafka

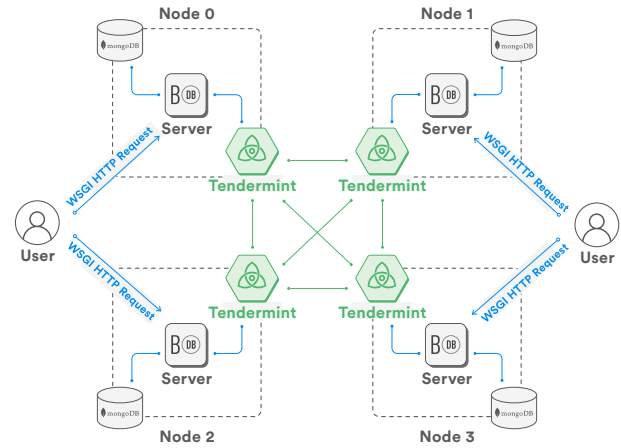


Figure 2: BigchainDB

recorded onto the blockchain. Hence, the correctness of the data can be verified based on the history stored in the blockchain.

An alternative design of Veritas [19], which is selected by us to be implemented in this paper, uses a shared verifiable table as its key storage. In this design, each node has a whole copy of the shared table and tamper-proof logs in the form of a distributed ledger, as shown in Figure 1. The ledger stores the update (write) logs of the shared table. Each node sends its local logs and receives remote logs via a broadcasting service. The verifiable table design of Veritas uses timestamp-based concurrency control. The timestamp of a transaction is used as the transaction log's sequence number, and each node has a watermark of the committed log sequence. When a transaction request is sent to a node, it first executes the transaction locally and buffers the result in memory. Then, it sends the transaction log to the other nodes via the broadcasting service. The node flushes the buffer of transactions and updates the watermark of committed logs as soon as it receives approval from all the other nodes.

BigchainDB [20] uses MongoDB [24] as its storage engine. That is, each node maintains its local MongoDB database, as shown in Figure 2. MongoDB is used due to its support for assets, which is the main data abstraction in BigchainDB. Tendermint [3] is used for consensus among the nodes in BigchainDB. Tendermint is a BFT consensus protocol and it guarantees that when one node is

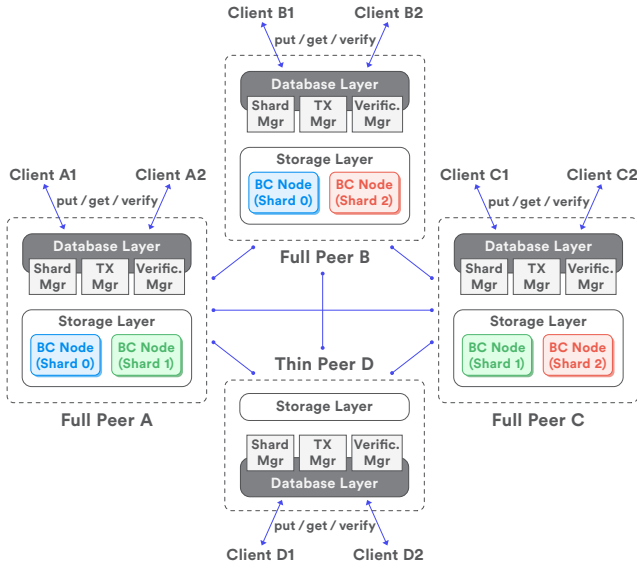


Figure 3: BlockchainDB

controlled by a malicious hacker, the MongoDB databases in the other nodes will not be affected. When a node receives an update request from a user, it first generates the results locally and makes a transaction proposal to be sent to the other nodes via Tendermint. The node commits the buffered result and responds to the user client as soon as most of the nodes in BigchainDB reach a consensus on this transaction.

BlockchainDB [16] adopts a design that builds a shared database on top of a blockchain. It is different from the other systems because it partitions the database into a few shards, as shown in Figure 3, such that the overall storage overhead is reduced. While some storage is saved, this design induces a higher latency since a data request may need to do an additional lookup to locate the corresponding shard. In BlockchainDB, each peer integrates a shard manager to locate the shard where a specific key is located. In terms of verification, it provides both synchronous (online) and asynchronous (offline) verification which is done in batches.

FalconDB [33] is a system that provides auditability and verifiability by requiring both the server nodes and the clients to keep a digest of the data. The server nodes of FalconDB keep the shared database and a blockchain to record the update logs of the shared database. Client nodes only hold the block headers of the blockchain kept by the server nodes. Using these headers, the clients are able to verify the correctness of the data obtained from the server nodes. These client nodes act as intermediaries between the users and the actual database.

ChainifyDB [38] proposes a new transaction processing model called Whatever-Ledger Consensus (WLC). Unlike other processing models, WLC makes no assumptions about the behavior of the local database. The main principle of WLC is to seek consensus on the effect of transactions, rather than the order of the transactions. When a ChainifyDB server receives a transaction from a client, it asks for help from the agreement server to validate the transaction and then sends a transaction proposal to the ordering server. The

ordering server batches the proposals into a block with FIFO order and distributes the block via Kafka [17]. When the transaction is approved by the consensus server, it will finally be executed in-order in the execution server's underlying database.

At a high level, Blockchain Relational Database [31] is very similar to Veritas [19]. However, in Blockchain Relational Database [31] (BRD), the consensus is used to order blocks of transactions, and not to serialize transactions within a single block. The transactions in a block of BRD are executed concurrently with Serializable Snapshot Isolation (SSI) on each node and they are validated and committed serially. PostgreSQL [12], which supports Serializable Snapshot Isolation, is used as the underlying storage engine in BRD. The transactions are executed independently on all the "untrusted" databases, but then they are committed in the same serializable order via the ordering service.

2.2 Ledger Databases

Different from hybrid blockchain database systems, ledger databases [1, 7, 26, 45, 48] are centralized, as the ledger is kept by a single organization. In this paper, we briefly describe some of the existing ledger databases. However, we are not evaluating and analyzing their performance.

Amazon Quantum Ledger Database [1] (QLDB) contains an immutable journal that documents every data change in a precise and sequential manner. The journal is made up of append-only blocks that are arranged in a hash chain. This means that data can only be appended to the journal and cannot be overwritten or deleted. The entire journal is designed as a Merkle Tree, allowing users to trace and check the integrity of data changes.

Immudb [7] is a lightweight, high-speed immutable database with built-in cryptographic proof and verification. Immudb is written in pure Go, with BadgerDB [13] as its storage engine. Badger is a fast key-value database implemented in pure Go, which uses an LSM tree structure. Immudb guarantees immutability by using a Merkle Tree structure internally where data is hashed with SHA-256. Moreover, immudb builds a consistency checker to check the correctness of data periodically.

Spitz [48] is a ledger database that supports a tamper-evident and immutable journal of transactions. It uses Forkbase [44] as its underlying storage, which provides a git-like multi-version control for data. Spitz provides a cell store for storing data and an append-only ledger to store the journal of transactions. Moreover, it builds a Merkle Tree based on the ledger to provide verifiability.

LedgerDB [45] is a centralized database from Alibaba Group. It uses TSA time notary anchors to provide auditability. These anchors are generated by a two-way peg protocol [43]. What is different in LedgerDB compared to the previous ledger databases is that it supports not only create, update, and query methods but also purge and occult methods for verifiable data. With these methods, LedgerDB aims to meet the requirements of the real world. However, it may destroy immutability while providing strong verifiability. As for the underlying storage, LedgerDB supports file systems including HDFS [40], key-values stores such as RocksDB [41], Merkle Patricia Tree [47] and a linear-structured append-only file system called L-Stream [45] which is specially designed for LedgerDB.

Table 1: Hybrid Blockchain Database Systems

| System | Storage | Consensus | User API |
|----------------------|----------------|----------------------|----------|
| Veritas (Kafka) [19] | Redis | Kafka (CFT) | KV-Store |
| Veritas (TM) | Redis | Tendermint (BFT) | KV-Store |
| BigchainDB [20] | MongoDB | Tendermint (BFT) | KV-Store |
| BlockchainDB [16] | RocksDB (Geth) | PoA Consensus (Geth) | KV-Store |
| FalconDB [33] | MySQL | Tendermint (BFT) | SQL |
| BRD [31] | PostgreSQL | BFT-SMaRt (BFT) | SQL |
| ChainifyDB [38] | MySQL | Kafka (CFT) | SQL |

2.3 Summary

In summary, hybrid blockchain databases are decentralized systems that have three key components, namely, (1) a shared database that uses a storage engine such as MySQL, MongoDB, Redis, among others, (2) a shared ledger replicated via a consensus mechanism such as Kafka (CFT) and Tendermint (BFT), among others, and (3) a user interface (API) that usually supports a simple key-value (KV) store with put and get operations. In Table 1, we summarize the technologies used for these three key components by the existing hybrid blockchain database systems.

3 SYSTEMS UNDER TEST

In this section, we describe the systems analyzed in this paper. We start with Veritas (Kafka) which uses Apache Kafka for inter-node communication. Then we describe our modification of this system into Veritas (TM) which uses Tendermint for the broadcasting service. Next, we present our implementation of BlockchainDB. We end by describing the existing implementation of BigchainDB.

3.1 Veritas (Kafka)

3.1.1 Overview. As shown in Figure 1, the key components of Veritas are the server nodes, client nodes, timestamp service, and broadcasting service. Our implementation, named **Veritas (Kafka)** uses Apache Kafka [17], a CFT pub-sub system, as broadcasting service. Next, we describe the functionality of each component of Veritas. **Server nodes** keep both the shared database and the full ledger that contains all the logs of the transactions done on the shared database. Server nodes handle query and update requests received via the client nodes. They also need to provide proofs of data integrity and accuracy as part of the verify requests. **Client nodes** act as intermediaries between users and server nodes. Upon getting a request, a client node first gets a global timestamp from the timestamp service. This timestamp represents the unique identifier of the transaction throughout its lifetime. **Timestamp service** generates global timestamps used by the client nodes. This service needs to provide unique monotonically increasing timestamps. Moreover, this service should exhibit strong availability and high performance. **Broadcasting service** gets local transaction logs from a server node and distributes them to all the other server nodes. Similar to the timestamp service, it should exhibit strong availability and high performance.

When a transaction contains an update operation, the server node first executes the transaction locally and buffers the result in memory. Then, it sends the transaction log to the broadcasting service which distributes it to the other server nodes. When the initial server node receives the approvals from all the other server

Table 2: Veritas User API

| Operation | API |
|-----------|--|
| Begin | BeginTxn(signature) -> txn_handler |
| Commit | Commit(sync / async) -> (txn_id, error_code) |
| Query | Query(txn_id) -> error_code |
| Set | Set(key, value) -> error_code |
| Get | Get(key) -> (value, version, error_code) |
| Verify | Verify(key) -> (root_digest, proof_path) |

nodes, it commits the buffered result to the local database and appends the transaction log to the ledger.

3.1.2 User API. Veritas treats all user requests as transactions [19]. Each transaction contains one or more operations. The operations supported by Veritas are shown in Table 2 and explained in the following lines. *Begin* starts a transaction using the user’s signature. It returns a transaction handler to deal with the next operations of the transaction. *Commit* finalizes a transaction. It has two modes of operation. The first mode is synchronous, where the user needs to wait for the transaction commit result. The second mode is asynchronous, where the user does not need to wait for the result of the Commit operation. When the user is using the asynchronous Commit, she needs to use the Query operation to get the result of the entire transaction. *Query* checks the status of a transaction. As all the transaction logs are stored in the ledger, Veritas can easily find whether the specified transaction is committed or aborted. *Set* updates the value of a specified key. The value is linked to the transaction unique identifier, which is recorded on the ledger. *Get* retrieves the value of a specified key. In our implementation of Veritas, we guarantee that the user reads the latest committed value of the key. *Verify* traces the proof path of the specified key in the Merkle Tree. Users can verify the correctness of the value by calculating the root digest of the proof path, and then compare the result with the expected root digest.

3.1.3 Implementation Details. We implemented Veritas (Kafka) in 804 lines of Go code, using Redis [5] (v6.2.1) as the underlying database for shared tables. Redis is an in-memory key-value store that exhibits very high throughput, of up to 100,000 operations per second [28]. To store the ledger, we use BadgerDB [13] which is a database engine based on LSM trees and optimized for SSD. BadgerDB stores the keys and values separately to reduce the IO cost. Moreover, it stores pointers to the values in the LSM tree to save the compaction cost.

In the original design of Veritas [19], the ledger is stored using a Merkle Tree [30]. While it is easy to verify that something is part of a Merkle Tree, it takes more effort to prove that something is not in the tree. For this reason, we adopt a Sparse Merkle Tree [10] in our implementation. This kind of tree contains a leaf for every possible result of a cryptographic hash function, making it easy to show that certain data is not part of the tree. On the other hand, a Sparse Merkle Tree needs more storage space compared to a Merkle Tree.

The timestamp service implementation is based on Timestamp Oracle (TSO) [32] which ensures that the clock assigned to an event will not repeat. In Veritas, this timestamp service ensures that a unique and monotonically increasing timestamp is applied

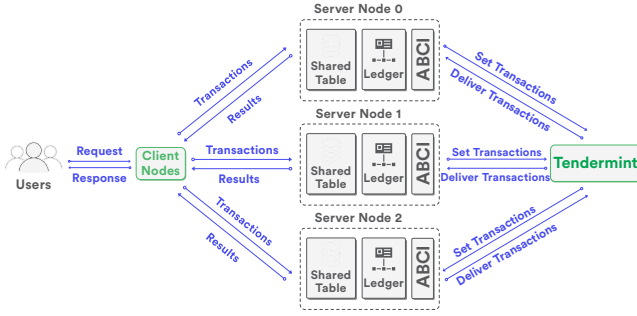


Figure 4: Veritas (TM) with Tendermint

to each transaction. Other systems, such as TiDB [23, 34], also use Timestamp Oracle. In our implementation, the TSO is crash fault-tolerant by using a persisted write-ahead log.

The broadcasting service is based on Apache Kafka [17] (v2.7.0), which is a distributed messaging system that provides high throughput. Kafka ensures crash fault tolerance (CFT) by persisting the messages to the disk and replicating them across nodes.

3.1.4 Complexity Analysis. Next, we analyze the complexity of Veritas (Kafka) in terms of the number of messages the nodes need to exchange per block of updates. We suppose there are N Veritas server nodes, K Kafka nodes, and the replication factor of Kafka is R . In practice, this replication factor is set to three [25]. In Veritas (Kafka), a node sends a block of updates to the Kafka service. Besides replicating it on R replicas, the Kafka service sends the block to $N - 1$ Veritas nodes. Each node checks and applies the updates, after which it sends an acknowledgement to Kafka. Then, Kafka sends each of the $N - 1$ acknowledgements to the other $N - 1$ nodes, resulting in a message complexity of $O(N^2)$ for Veritas (Kafka). This result is backed by experimental evaluation in Section 4.3.

In terms of storage complexity, we note that each block is persisted in the Sparse Merkle Tree of each Veritas node, and also on R Kafka nodes. Hence, the storage complexity is $O(N + R)$.

3.2 Veritas (TM)

3.2.1 Overview. In general, blockchains operate in a Byzantine environment [14, 36]. Hence, using a CFT broadcasting service or consensus protocol, such as Apache Kafka or Raft, is unsuitable. That is why we implement Veritas with Tendermint [3], a BFT consensus protocol used by other systems, such as BigchainDB [20] and FalconDB [33]. Tendermint is a BFT middleware that supports a state transition machine. Similar to other BFT systems, Tendermint supports up to $1/3$ malicious nodes in the network.

Figure 4 depicts the architecture of **Veritas (TM)**. We use Tendermint Core as the consensus and broadcasting service and an Application BlockChain Interface (ABCI) application integrated with the Veritas node to interact with Tendermint. A *Get* request is served directly by the Veritas node. In contrast, a *Set* transaction is sent to Tendermint via ABCI. When the transaction is included in a block and delivered back to the Veritas node via ABCI, it is applied to the local (Redis) database. This mechanism ensures a serial concurrency control in Veritas (TM). Different from Veritas (Kafka), we do not use a timestamp oracle (TSO) service in Veritas (TM). Such a

TSO service is centralized in nature, thus, incompatible with a BFT setting.

3.2.2 Implementation Details. Veritas (TM) is implemented in 517 lines of Go code, with Tendermint [3] (v0.35.0) as its consensus service, which is to be consistent with BigchainDB. The ledger module, verifiable data structure, concurrency control, and local database are the same as for Veritas (Kafka). That is, we use a Sparse Merkle Tree as the verifiable data structure and Redis (v6.2.1) as local database.

3.2.3 Complexity Analysis. The message complexity of Veritas (TM) is determined by the complexity of Tendermint, which is shown to be $O(N^3)$ by a recent study [2], where N is the number of nodes in the system. We shall see in Section 4.3 that the performance of Veritas (TM) decreases drastically with the number of nodes. We attribute this to the message complexity of Tendermint. The storage complexity of Veritas (TM) is $O(N)$ since each block is stored on all N nodes.

3.3 BlockchainDB

3.3.1 Overview. As shown in Figure 3, a BlockchainDB [16] node is mainly composed of a storage layer and a database layer. The storage layer uses an off-the-shelf blockchain as a shared database which provides tamper-proof and de-centralized storage. A blockchain connector component, which is implemented for different blockchain clients, handles the interactions with the blockchain network. The database layer is built on top of the storage layer with a simple read/write interface to access data from a shared table. The database layer in BlockchainDB is responsible for handling the *Set* and *Get* requests from clients. Different from all the other systems under test, BlockchainDB supports sharding. A sharding manager component defines a partition scheme based on a hash algorithm and stores the connection information for each shard.

3.3.2 User API. As shown in Table 3, BlockchainDB [16] supports three operations, namely, *Set*, *Get*, and *Verify*. The *Set* method writes a key-value pair to the given shared table and returns a corresponding blockchain transaction id. The *Get* method retrieves the data corresponding to a key from the given shared table. The *Verify* method allows users to check the status of the given transaction id to verify if the operation was successfully committed.

3.3.3 Implementation Details. In this paper, BlockchainDB is implemented in Go, where each node runs an RPC server serving clients requests. We use a private Ethereum (geth/v1.8.23-stable) blockchain with Proof-of-Authority (PoA) consensus as the storage for experiments. A KVStore contract as defined in [16] is installed on the Ethereum network.

3.3.4 Complexity Analysis. BlockchainDB is a sharded system, hence, we analyze the communication and storage overhead for a single shard. The total overhead is the sum of the overheads of each shard. Given the Ethereum backend used by BlockchainDB, the communication complexity is $O(N)$, where N is the number of nodes in a shard. This is due to the block broadcasting phase in Ethereum. The storage complexity is $O(N)$ since the ledger is stored by each node.

Table 3: BlockchainDB User API

| Operation | API |
|-----------|--------------------------------|
| Set | Set(table, key, value) -> void |
| Get | Get(table, key) -> value |
| Verify | Verify() -> bool |

Table 4: BigchainDB User API

| Operation | HTTP Method | URL |
|-----------|-------------|---------------------------------------|
| Query | GET | /api/v1/transactions/{transaction_id} |
| Create | POST | /api/v1/transactions?mode={mode} |
| Transfer | POST | /api/v1/transactions?mode={mode} |

Table 5: A Comparison of Features in the Systems Under Test.

| | Veritas (Kafka) | Veritas (Tendermint) | BlockchainDB | BigchainDB |
|------------------------|--------------------|----------------------|----------------------|-----------------------|
| <i>Fault Tolerance</i> | CFT (Kafka) | BFT (Tendermint) | BFT (via Ethereum) | BFT (Tendermint) |
| <i>Immutability</i> | Append-only Ledger | Append-only Ledger | Ledger (Ethereum) | Ledger (Tendermint) |
| <i>Verification</i> | Sparse Merkle Tree | Sparse Merkle Tree | Merkle Patricia Trie | Merkle Tree |
| <i>Concurrency</i> | MVCC | Serial | Serial | Blockchain Pipelining |
| <i>Storage</i> | NoSQL (Redis) | NoSQL (Redis) | Ethereum (geth) | NoSQL (MongoDB) |
| <i>Replication</i> | Full Replication | Full Replication | Sharding | Full Replication |

3.4 BigchainDB

3.4.1 Overview. As shown in Figure 2, a BigchainDB node consists of three key components, namely, the server, consensus component (Tendermint), and local database (MongoDB). **BigchainDB Server** provides an HTTP service that handles user requests. This service uses the Flask [21] web application framework working with Web Server Gateway Interface (WSGI) of Unicorn [6] to expose an HTTP API. **Tendermint Consensus Component** provides a broadcasting service for the transaction blocks. It has the role of a bridge among BigchainDB nodes and it is responsible for proposing new transaction blocks and ensuring that all nodes agree on a block in a Byzantine fault-tolerant manner. After validating a transaction block, the Tendermint component sends a commit message to the local BigchainDB server to signal the commit of the transactions. **MongoDB Database Component** persists the data which can only be modified by the local BigchainDB server. The exposed API of the MongoDB component is determined by the local BigchainDB server. Different MongoDB instances may provide different functions.

BigchainDB uses a concept called blockchain pipelining, which improves scalability when voting for the next blocks. In a blockchain, the transaction blocks are ordered which means that nodes cannot vote for a new block while the current block is undecided. This is because the new block needs to have a reference to a decided block. In BigchainDB, the blocks are also ordered, but server nodes are allowed to vote for a new block even if the current block is undecided. Using a voting list created at the same time with a block, expected voters are tracked while the list contains a reference to the current unsettled block. However, when voting for a block with an undecided parent block, a node has to verify that the block does not contain transactions with dependencies in the undecided block. This is a form of concurrency control adopted by BigchainDB [20]. We shall see in the next section that the validation process in BigchainDB slows down the entire system.

The transaction flow in BigchainDB can be described as follows. When a BigchainDB server receives an HTTP request with a transaction, it first performs some checks to validate the transaction. That is, the node checks if this is not a duplicate transaction in both

the transaction queue and the ledger. It also checks the inputs and outputs of the transaction. Next, it calls the *Broadcast API* provided by the local Tendermint component for the broadcasting of transactions. After receiving the commit message of the transaction from the Tendermint component, the server updates the state of the local MongoDB instance.

3.4.2 User API. As shown in Table 4, BigchainDB supports three operations: query, create, and transfer. *Query* retrieves the details of a transaction based on its transaction id. If the transaction has been committed, a BigchainDB server returns the details of the transaction. Otherwise, a 404 error code is returned. *Create* has the role to create assets for the specified user and to store them in BigchainDB. A Create transaction supports three modes, namely, async, sync, and commit. The default mode is async where a BigchainDB server responds before the Tendermint component validates the transaction. In sync mode, a BigchainDB server responds after the transaction block has been committed and the state in the MongoDB instance has been modified. Lastly, in commit mode, a BigchainDB server responds after checking the validation process of the block containing the transaction. Finally, *Transfer* has the role to transfer assets from one user to another. It also provides the same three modes for transaction processing as the Create API.

3.4.3 Implementation Details. In this paper, we use the open-source BigchainDB [20] (v2.2.2) with MongoDB [24] (v4.4.4) and Tendermint [3] (0.31.5). In addition to the standard system, we also evaluate BigchainDB with Parallel Validation feature (**BigchainDB (PV)**) which adopts an optimistic approach for validating transactions. This approach groups the transactions based on the assets they are touching and supposes there are no conflicts across these groups. Hence, transaction validation in groups can be parallelized on multiple CPU cores, becoming faster.

3.4.4 Complexity Analysis. The message complexity of BigchainDB is the same as that of Tendermint, which is $O(N^3)$ [2], where N is the number of nodes. The storage complexity is $O(N)$ because the ledger is stored on each node.

3.5 Summary

The design and implementation of a hybrid blockchain database system require choices for the consensus/communication mechanism, ledger data structure, persistent storage database, concurrency control, and replication. In Table 5, we compare the design choices in our Veritas implementations, BlockchainDB, and BigchainDB.

In BlockchainDB, the use of Ethereum blockchain as the underlying database has implications on immutability and verification since the ledger of Ethereum uses a Merkle Patricia Trie. In contrast, our Veritas implementations use a custom Sparse Merkle Tree for the ledger, which is integrated with the server node. For concurrency control, the systems under test use different approaches. Veritas (TM) and BlockchainDB use serial commit, while BigchainDB optimizes the serial commit with blockchain pipelining. Veritas (Kafka) adopts MVCC. At the storage level, both Veritas flavors use Redis, BlockchainDB uses Ethereum, while BigchainDB uses MongoDB. All systems, except BlockchainDB, use full replication (no sharding). Even for BlockchainDB, we use full replication in our experiments.

4 PERFORMANCE ANALYSIS

In this section, we analyze the performance of the five systems described in the previous section, namely, Veritas (Kafka), Veritas (TM), BigchainDB, BigchainDB (PV), and BlockchainDB. Next, we describe the experimental setup.

4.1 Setup

All the experiments are executed on a local machine under Ubuntu 18.04 operating system (OS). The machine has 256 physical CPU cores, 8 TB of RAM, and 3TB of hard-disk (HDD) storage. Using iostat tool, the IOPS of the machine is estimated at 5.35. All the server and client nodes of the systems under test are running on this machine in Docker containers on different CPU cores.

To evaluate the performance of the five systems under test, we send 100,000 transactions to each system. We send multiple transactions in parallel to a system, based on a concurrency parameter that represents the number of clients. The transactions are evenly distributed to different server nodes in the system. To compute the throughput, we record the start time of the first transaction and the completion time of all the transactions. Then, we record the number of successfully committed transactions and compute the throughput by dividing this number by the previously recorded time interval. Note that we only consider successful transactions when computing the throughput, but there may be failures as well. We repeat each experiment three times and report the average. Before executing the transactions, we load 100,000 key-value records of 1,000 bytes each into each system.

We use the Yahoo! Cloud Serving Benchmark [8] (YCSB) dataset which is widely used to benchmark databases and blockchains [15, 36]. YCSB supports common database operations such as write (insert, modify, and delete), and read. While YCSB defines six workloads, in our experiments, we choose three of these six workloads. These three workloads are (i) *Workload A* which consists of 50% update operations and 50% read operations, (ii) *Workload B* which consists of 5% update operations and 95% read operations, and (ii) *Workload C* which consists of 100% read operations. Moreover, we

use three key distributions, namely, (i) *uniform distribution* which operates uniformly on all the keys, (ii) *zipfian distribution* which operates frequently only on a sub-set of the keys, and (iii) *latest distribution* which operates on the latest used keys. We use Workload A and uniform distribution by default, unless otherwise specified.

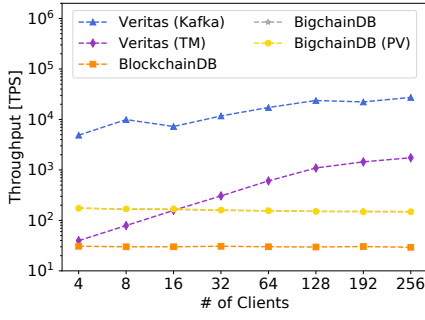
The benchmarking tools are implemented by us in Go using *goroutines* [9] and a *channel* [9] as a concurrent safe request queue. A channel allows goroutines to synchronize without explicit locks or condition variables. Each benchmarking client is represented by a goroutine and it gets a new request from the channel once it completes the current request.

4.2 Effect of Number of Clients

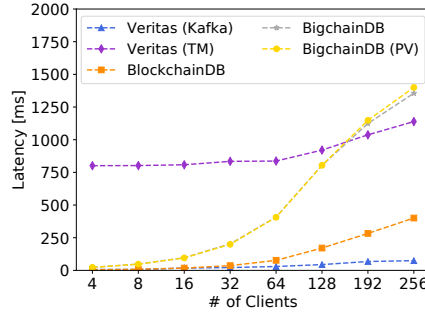
We first analyze the effect of increasing the number of client nodes to determine the saturation points of the systems. In this experiment, we set the number of server nodes in each system to four. This number is derived based on the fact that BFT systems supporting up to f faulty nodes need to have a total of at least $3f + 1$ nodes. Hence, for tolerating one faulty node, we need four nodes in the system. For consistency, we also set the number of Veritas (Kafka) server nodes to four. We use YCSB Workload A with uniform distribution and set the block size of the ledger to 100 transactions. We then increase the number of clients from 4 to 256.

Figures 5a and 5b show the effect of increasing the number of clients on throughput and latency, respectively. We observe that the throughput of the BFT systems plateaus after using a certain number of clients, while for Veritas (Kafka) it is growing even if at a slow pace. These results are partially correlated with the latency: we observe a sharper increase in the latency of the BFT systems when using more than 32 clients. Compared to BigchainDB, the throughput of Veritas (TM) becomes much higher starting from 32 clients. Even if both Veritas (TM) and BigchainDB use Tendermint as the underlying consensus protocol, their performance is very different. Veritas (TM) achieves its top performance of 1,742 TPS with 256 clients, while BigchainDB only achieves 175 TPS and this when using four clients. Increasing the number of clients in BigchainDB results in a sharp increase in latency while the throughput remains relatively constant.

To investigate the reasons for BigchainDB's low performance, we analyze the time spent in each major component of a hybrid blockchain database system. As presented in Section 2, each node of such a hybrid system has an underlying shared *database*, a *ledger* data structure and storage, and a *consensus* or broadcasting component. As shown in Figure 6, Veritas (Kafka), Veritas (TM), and BlockchainDB spend 45%, 55%, and 99% of their time in the consensus component. On the other hand, BigchainDB spends 38% of its time validating transactions, as part of the ledger component. For example, BigchainDB checks for duplicate transactions both in the transaction queue (in memory) and in the database (in MongoDB, on disk). This operation is very costly and implies sequential processing. That is why the latency increases significantly when more transactions are sent to the system at the same time (when increasing the number of clients). In BigchainDB (PV), transactions are checked in parallel but the overhead of parallel processing is around 11%. This, together with other ledger operations lead to 33% of the time spent in the ledger component. The database part in



(a) Throughput



(b) Latency

Figure 5: The Effect of the Number of Clients

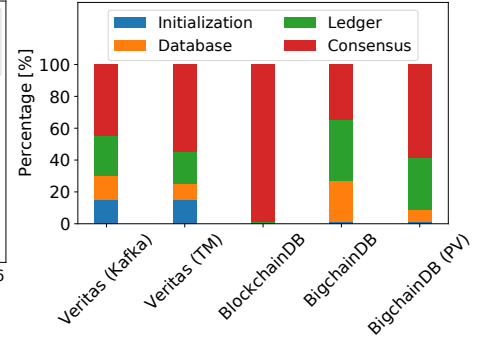
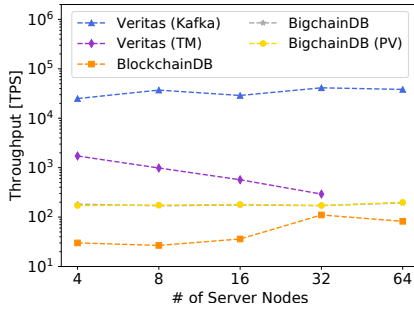
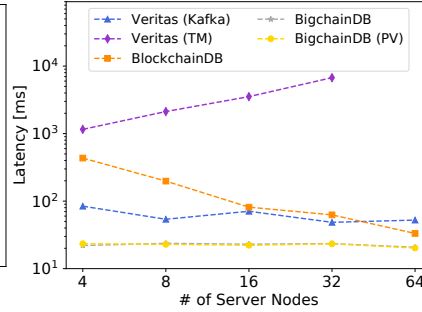


Figure 6: Breakdown per Component



(a) Throughput



(b) Latency

Figure 7: The Effect of the Number of Server Nodes

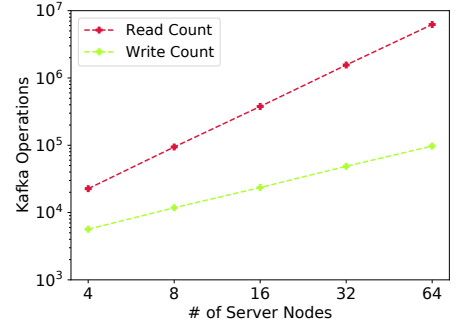


Figure 8: The Effect of Veritas Server Count on Kafka Operations

BigchainDB (PV) is much lower compared to BigchainDB due to the bulk storage of transactions. However, the consensus accounts for 58.5% in BigchainDB (PV) compared to 35% in BigchainDB. We attribute this to the larger consensus message size in BigchainDB (PV) compared to BigchainDB.

Secondly, we observe that the performance of Veritas (Kafka) is more than 10 \times higher compared to the BFT systems. In particular, Veritas (Kafka) exhibits a maximum of 27,335 transactions per second (TPS) when 256 clients are used, while Veritas (TM) achieves 1,742 TPS with 256 clients. BigchainDB and BigchainDB (PV) exhibit a maximum of 175 and 174 TPS, respectively, when using four clients. On the other hand, BlockchainDB achieves only 30 TPS, but this is expected due to the use of Ethereum as the underlying data storage. Even when using the PoA consensus, the throughput of Ethereum is below 100 TPS [15]. In terms of average latency, the corresponding values are 74, 1139, 400, 23, and 23 milliseconds (ms) for Veritas (Kafka) with 256 clients, Veritas (TM) with 256 clients, BlockchainDB with 256 clients, BigchainDB and BigchainDB (PV) with four clients, respectively.

Thirdly, we observe that the throughput of Veritas (TM) is 10 \times higher compared to BigchainDB. As shown in Table 5, the difference between these two systems is at storage and concurrency control layers. At the storage layer, Veritas (TM) is using Redis, while BigchainDB is using MongoDB. Redis has a higher throughput than MongoDB, being an in-memory database. However, the throughputs of Veritas (TM) and BigchainDB are far lower than

what Redis and MongoDB can support. To investigate if MongoDB has a significant negative impact on the performance, we replace Redis with MongoDB in Veritas (TM) and run the benchmarking again. The performance of Veritas (TM) with MongoDB is up to 22% lower compared to using Redis. In conclusion, the huge difference between Veritas (TM) and BigchainDB is due to the design and implementation of the latter, especially the complex and redundant transaction validation mechanism.

4.3 Effect of Number of Server Nodes

Next, we analyze the effect of increasing the number of server nodes on throughput and latency. We start from four nodes, the minimum number of nodes in a BFT system to support one faulty node, and increase the number of nodes up to 64. A network of 64 nodes can support 21 faulty nodes in a BFT system. Similar to the previous experiment, we use YCSB Workload A with uniform distribution and a block size of 100 transactions. We set the number of clients to 256 for Veritas (Kafka), Veritas (TM), and BlockchainDB, and to four for BigchainDB and BigchainDB (PV).

As shown in Figure 7, we observe the same gap in performance among Veritas (Kafka) and the BFT systems. Note that the latency of BlockchainDB only includes the time it takes for a request to be included in the transaction pool of Ethereum. It does not include the block confirmation time. The throughput of Veritas (TM) drops from 1,722 TPS on four nodes to 292 TPS on 32 nodes. On 64 nodes, Tendermint (v0.35.0) stops working due to synchronization errors

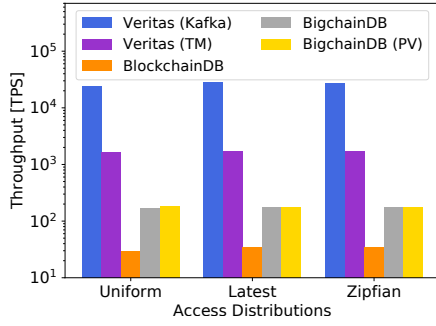


Figure 9: The Effect of Key Access Distribution on Throughput

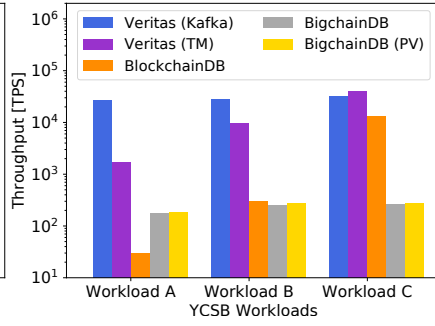


Figure 10: Performance of Different YCSB Workloads

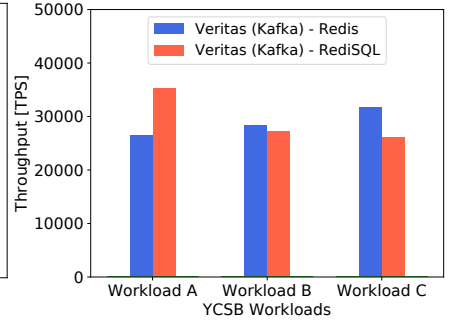


Figure 11: The Effect of the Underlying Database

among the peers. Again, this trend can be explained by the high message complexity of Tendermint, namely, $O(N^3)$, where N is the number of nodes. In BigchainDB and BigchainDB (PV), the impact of Tendermint is amortized because the system is slow in sending new transactions to the consensus component. As explained before, BigchainDB and BigchainDB (PV) perform complex and redundant transaction validations which slow down the entire system.

In Veritas (Kafka), the throughput increases with the number of server nodes, from 25,009 TPS on four nodes to 41,262 TPS on 32 nodes. On 64 nodes, the throughput decreases to 38,283 TPS. We attribute this to the interplay between the number of server nodes and the Kafka broadcasting service. With four server nodes, there are not enough transactions exchanges to saturate Kafka which supports up to 200,000 messages per second. Also, there are too few server nodes to process the transactions as opposed to 32 nodes that can process more transactions. But when we increase the number of nodes beyond 32, Kafka starts to become a bottleneck. For example, with 64 nodes, there are close to 100,000 writes and 6,200,000 reads in the Kafka service.

In Figure 8, we plot the number of read and write operations in Kafka as the number of server nodes in Veritas (Kafka) grows, using a logarithmic scale. The number of write operations grows linearly with the number of server nodes because, for each transaction sent by a node, the other $N - 1$ nodes need to write (send) an acknowledgment to Kafka. On the other hand, the number of read operations grows quadratically with the number of server nodes. This is because each of those $N - 1$ acknowledgments are read by all the other $N - 1$ nodes in the system. This results in a message complexity of $O(N^2)$ in Veritas (Kafka), where N is the number of server nodes.

4.4 Effect of Access Distributions

In this section, we evaluate the performance of the five systems with different access distributions, namely uniform, latest, and zipfian with a coefficient of 0.5. As shown in Figure 9, we do not observe significant differences among these distributions. This can be explained by the fact that there is no significant number of aborted transactions in the evaluated systems. For example, Veritas (TM) and BlockchainDB exhibit no aborted transactions due to the serial transaction commit. Even in Veritas (Kafka) with MVCC, there are up to 20 aborted transactions out of 100,000 total transactions,

when using the zipfian distribution. We attribute this low number of conflicts to the fact that the version control is based on the centralized global timestamp service. This, combined with a fast Kafka delivery service, make conflicts unlikely.

In Veritas (Kafka), serialization is guaranteed by the broadcasting service. Server nodes send both new blocks and approvals of blocks concurrently, with no blocking. When a server node receives a new block, it first checks the conflicts between the transactions in that block and the current state of the underlying database. If there is no conflict, it sends an approval message with its signature via Kafka and buffers this block in memory. If there are conflicts, the server node sends an abort message via Kafka. The block is finally committed once the server node receives all the approval messages from the other server nodes. Hence, conflict checking and new block generation happen in parallel in Veritas (Kafka). But this leads to more conflicts and these conflicts are detecting during the checking phase.

4.5 Effect of YCSB Workloads

Next, we evaluate the performance of the five systems with three different YCSB workloads, namely, Workload A, Workload B, and Workload C. The number of transactions in each block of the shared log is 100, the number of server nodes is set to four, and the number of clients in each system are those that exhibit the best performance.

As expected, all the systems exhibit higher performance when the number of read operations is higher. Recall that the proportion of read operations is 50%, 95%, and 100% in Workloads A, B, and C, respectively. As shown in Figure 10, the gap between Veritas (Kafka) and Veritas (TM) becomes smaller as the number of read (or *Get*) operations is higher. For example, Veritas (TM) achieves higher throughput with Workload C compared to Veritas (Kafka). This is because *Get* operations are served immediately by Veritas (TM), while Veritas (Kafka) needs to apply the MVCC mechanism that introduces a delay. For the same reason of serving *Get* operations immediately, BlockchainDB achieves an impressive throughput of 13,124 TPS with Workload C, compared to just 30 TPS with Workload A. On the other hand, BigchainDB achieves only 257 TPS with Workload C. We attribute this to the complex asset query in MongoDB used by BigchainDB to serve a *Get* operation.

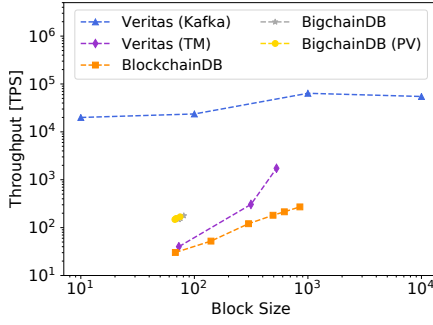


Figure 12: The Effect of Block Size on Throughput

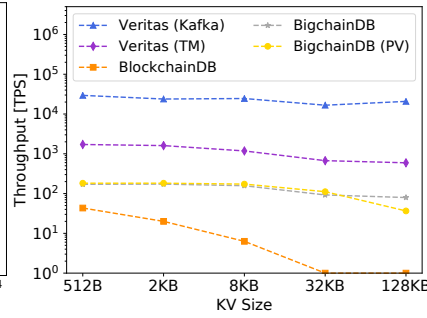


Figure 13: The Effect of Record Size on Throughput

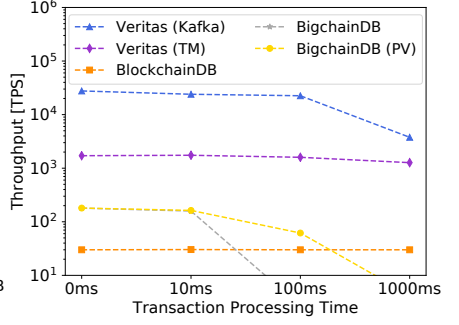


Figure 14: The Effect of Processing Time on Throughput

4.6 Effect of the Underlying Database

The developers of hybrid blockchain database systems have the choice to use either a relational database or a simple key-value store as the underlying database of the system. In our implementations of Veritas, we choose Redis, a fast in-memory key-value store. Here, we compare Veritas (Kafka) using Redis with Veritas (Kafka) using an SQL database. For fairness, we choose an in-memory SQL engine, namely RediSQL [35]. RediSQL is a fast, in-memory SQL database based on Redis, which supports standard SQL semantics.

As shown in Figure 11 for three YCSB workloads, the throughput of Veritas (Kafka) using RediSQL decreases with a higher proportion of read operations. Even with an index on the key in RediSQL, the latency of read (Get) is high, and this leads to lower throughput when there are more read operations in the workload. On the other hand, it is interesting to observe that Veritas (Kafka) using RediSQL has higher throughput compared to Veritas (Kafka) using Redis on Workload A. This is an interesting effect of the interplay between fast read operations and the fine locking used in Veritas (Kafka). Read operations are fast in Redis and this leads to more lock conflicts in write (Set) operations. On the other hand, read operations are slower in RediSQL compared to Redis and there are fewer lock conflicts in Veritas (Kafka). In conclusion, using an SQL or NoSQL database as the underlying storage may have little to moderate effect on the system performance because the consensus module is always the bottleneck. However, when SQL semantics are needed, the developer can use an SQL database rather than a NoSQL database.

4.7 Effect of Block Size

In the existing blockchain literature [37, 39], block size is shown to have a significant impact on the performance of a blockchain system. Thus, we evaluate this impact on our hybrid blockchain database systems by increasing the block size while using YCSB Workload A. We note that in Tendermint and Ethereum it is not possible to fix the block size because it fluctuates with other parameters, such as the request rate in Tendermint and the gas limit per block in Ethereum. For the systems using these underlying consensus, we present a few samples that exhibit different average block size.

As shown in Figure 12, the throughput of all the systems increases when the block size increases as there are more transactions sent per each communication round among the server nodes. For example, the throughput of Veritas (TM) increases from 40 TPS

with 73 transactions per block, to 1.734 TPS with 530 transactions per block. In Veritas (Kafka), the throughput grows from 20,000 TPS with 10 transactions per block to 64,000 TPS with 1,000 transactions per block. Then, the throughput decreases to around 55,000 TPS with 10,000 transactions per block. We attribute this decrease to a slightly higher average transaction latency when there are more transactions per block.

A bigger block size leads to higher average transaction latency. This is because when packing more transactions in a block, these transactions need to wait more before being validated, hence, their processing latency is higher. We also observe that the latency gap between the CFT- and BFT-based systems becomes wider when increasing the block size, while the throughput gap becomes smaller. Therefore, the user needs to make a trade-off between throughput and latency when increasing the block size. This trade-off depends on the application scenario.

4.8 Effect of Transaction Characteristics

In this section, we analyze the effect of two key transaction characteristics on the performance of hybrid blockchain database systems. These two characteristics are the transaction size and the transaction processing time, respectively. First, we analyze the effect of increasing the transaction size in terms of the total key-value size from 512B to 128kB. Recall that our previous experiments are done with a key-value size of 1kB.

Indeed, the throughput of all the systems decreases when the key-value size increases. Moreover, there is a sharper decrease in throughput for key-value sizes greater than 8kB. This is mainly due to the networking overhead of bigger records. We note that some of the consensus frameworks used in this paper, namely Apache Kafka and Tendermint, have a recommended maximum message size of 1MB [25, 42]. Even more, the recommended message size in Kafka is 1kB [27]. But a consensus message includes a block of transactions, hence, the larger the key-value size, the lower the number of transactions per block. This leads to higher transaction propagation delay and, in the end, the throughput drops and the average transaction latency increases. BlockchainDB stops working when the key-value size is greater or equal to 32kB. This is because the gas limit of 10,000,000 is not enough to process transactions with such large sizes. This is an important limitation of using Ethereum as the underlying blockchain in BlockchainDB.

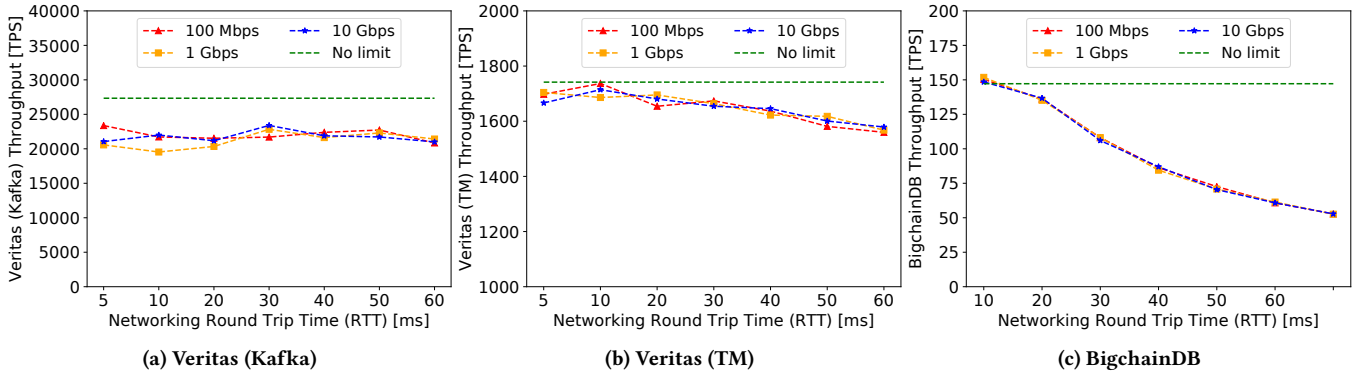


Figure 15: The Effect of Networking on Throughput

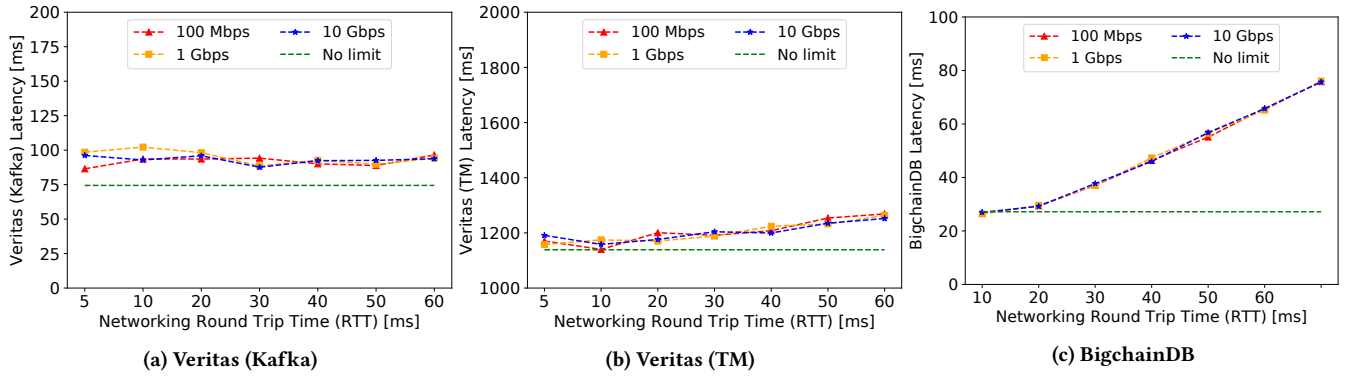


Figure 16: The Effect of Networking on Latency

Second, we analyze the effect of compute-heavy transactions by simulating heavy processing in the existing Set operations supported by the existing systems. This is done using a delay of up to one second in each transaction. While compute-heavy transactions are more common in blockchains' smart contracts, simple key-value stores may also implement complex key-value processing and validation which may increase the processing time.

As expected, the throughput of the systems decreases significantly when the transactions take longer to process. However, in Veritas (TM) the impact is not so obvious as in the other systems. For example, the throughput drops from around 1,700 TPS to 1,250 TPS when the transaction processing time increases from zero to one second in Veritas (TM). This could be explained by the fact that the average latency of Veritas (TM) is already higher than one second, and this is due to the consensus. Adding more execution time to the *Set* requests does not affect the consensus which partially overlaps with transaction processing. Hence, the increase in the average transaction latency is only 400 ms when the transaction execution time is one second in Veritas (TM).

4.9 Effect of Networking

In the previous sets of experiments, we assumed ideal networking conditions among all the nodes in the hybrid blockchain database systems. However, this is not the case in real-world deployments,

especially in cross-country blockchain systems. Hence, in this section, we evaluate the effect of different networking conditions on the performance of the five systems under test. For this, we use Docker containers and Open vSwitch (OVS) to change the networking bandwidth and latency among the nodes in the systems. In particular, we limit the bandwidth to 100 Megabits per second (Mbps), 1 Gigabit per second (Gbps), and 10 Gbps. At the same time, we increase the networking delay of a node from 2.5 to 30 ms, resulting in a Round Trip Time (RTT) in the interval 5 to 60 ms.

The throughput and latency results are shown respectively in Figures 15 and 16 for Veritas (Kafka), Veritas (TM), and BigchainDB. The trends for BlockchainDB and BigchainDB (PV) are not shown due to lack of space, but they are similar to those of Veritas (TM) and BigchainDB, respectively. In all the systems, adding delay and limiting the bandwidth leads to lower throughput and higher transaction latency. For all the systems, except Veritas (Kafka), we observe that transaction latency increases with growing RTT. This is expected since delay is introduced among all the messages exchanged by the peers. This higher latency has some impact on throughput, especially in Veritas (TM) and BigchainDB. In Veritas (TM), the throughput drops from 1,742 TPS to around 1,550 TPS when the RTT is 60 ms. In BigchainDB, the throughput drops from 150 TPS to around 50 TPS when the RTT is 60 ms. Otherwise, hybrid blockchain database systems can function well even in 100 Mbps networks.

4.10 Summary

In this section, we have shown that our implementation of Veritas (Kafka) can achieve close to 30,000 TPS with a block size of 100 and close to 65,000 TPS with a block size of 1,000 transactions, both with sub-second latency. On the other hand, the systems using Tendermint have low throughput, of up to 1,750 TPS, and high latency of up to 2 seconds. Our implementation of Veritas (TM) exhibits higher throughput than BigchainDB and BigchainDB (PV). This is because BigchainDB has a complex and redundant transaction validation mechanism which slows down the entire system. BlockchainDB has a throughput of less than 100 TPS because it uses Ethereum as its underlying storage engine. In summary, CFT systems, which are closer to the design of distributed databases, achieve higher performance compared to the systems based on BFT consensus but do not offer protection in Byzantine environments.

5 DISCUSSION

In this section, we discuss some trade-offs that users and developers need to consider when using and designing a hybrid blockchain database system. We look at these trade-offs from the perspective of performance, security, and ease of development.

CFT vs. BFT. Adopting a CFT or BFT consensus protocol depends on the environment where the system is supposed to operate. Nonetheless, it is well-known that BFT consensus is much slower than CFT consensus [11, 36]. In this paper, we show that the gap between the CFT and BFT systems throughput is 5 – 500×. We also show that the performance of both CFT and BFT systems degrades under poor real-world networking conditions. While typical blockchains employ BFT consensus protocols due to their security guarantees in a Byzantine environment, some permissioned blockchains such as Hyperledger Fabric [18] and R3 Corda [22] are employing CFT consensus protocols, while delegating the security aspects to the blockchain access/membership layer.

When choosing BFT, the main focus of a developer is on optimizing the performance since the security guarantees are relatively high. Query optimization and using high-performance underlying storage are the common methods of optimizing such systems. Optimizing the performance of the BFT consensus is an important aspect that is studied in recent works [3, 29, 46]. In contrast, CFT protocols usually provide high performance and developers need to add verifiable components to provide security guarantees such as building an append-only ledger to provide verifiability.

In summary, CFT hybrid blockchain database systems should be used when there is a demand for high performance without strict security guarantees, such as in online transaction processing (OLTP) systems. When starting with an existing CFT system, the developer can add blockchain features, such as an append-only ledger, to improve security. BFT hybrid blockchain database systems should be used when there is a need for high security, such as for digital assets storage. When starting with an existing BFT system, the developer can add database features, such as parallel validation and query optimization, to improve the performance.

SQL vs. NoSQL. When designing a hybrid blockchain database system, the developer needs to consider the data access model from both the user-level API and underlying storage perspectives. For example, the relational model and its de facto query language,

SQL, are widely used in the industry. But adding an SQL API in a hybrid blockchain database system is tedious because we need a layer to parse the SQL statements to basic operations and an optimized execution plan to get higher performance. BRD [31] uses this approach to support SQL. Another option is to add SQL support only to the underlying database, while the SQL queries are forwarded to this database. For example, FalconDB [33] and ChainifyDB [38] use MySQL to provide such a relational model in the underlying database.

On the other hand, many Internet-scale systems use simpler data models, such as key-value stores. For example, NoSQL databases support simpler data models while offering flexibility and scalability. However, they do not offer ACID properties. In this paper, we analyze Veritas and BigchainDB, two systems that rely on NoSQL systems. However, we show that the underlying database, being it NoSQL Redis or SQL RedisSQL, has little impact on the performance of Veritas (Kafka) and Veritas (TM).

Using SQL or NoSQL is a trade-off between integrity and flexibility. In SQL, the user must define the table and field structures before adding data. For example, the primary key, index, trigger, and stored procedure of the table must be defined to provide data constraints and prevent the database from storing invalid data. The table structure can be updated after being defined, but this is complicated if there are major structural changes. In NoSQL, data can be added at any time and anywhere, without the need to define the tables first. This provides more flexibility but may result in storing invalid data. In addition, SQL allows the integrity of the database to be constrained by defining foreign keys while there are no options for data integrity constraints in NoSQL.

6 CONCLUSIONS

In this paper, we survey the existing hybrid blockchain database systems proposed by the database community in the last few years. We then conduct a detailed performance analysis of five implementations of hybrid blockchain database systems based on Veritas [19], BlockchainDB [16], and BigchainDB [20]. Our flexible implementation of Veritas allows us to change the key components of the system, such as the consensus middleware and the underlying database, in order to analyze their effect on performance. For example, we replace the CFT Apache Kafka middleware with the BFT Tendermint and observe that the performance drops by more than 15×. However, Tendermint offers security guarantees in a Byzantine environment. On the other hand, we replace the underlying Redis database with RedisSQL and MongoDB and observe moderate performance differences. This, together with a breakdown of the time spent in each major component of a hybrid blockchain database system, shows that the consensus protocols are the source of poor performance in these systems.

ACKNOWLEDGMENTS

This research/project is supported by the National Research Foundation, Singapore under its Emerging Areas Research Projects (EARP) Funding Initiative. Any opinions, findings and conclusions or recommendations expressed in this material are those of the author(s) and do not reflect the views of National Research Foundation, Singapore.

REFERENCES

- [1] Amazon. 2021. Amazon Quantum Ledger Database (QLDB). <https://aws.amazon.com/qldb/>.
- [2] Yackolley Amoussou-Guenou, Antonella Del Pozzo, Maria Potop-Butucaru, and Sara Tucci Piergiovanni. 2019. Dissecting Tendermint. In *Proc. of 7th International Conference on Networked Systems (Lecture Notes in Computer Science)*, Mohamed Faouzi Atig and Alexander A. Schwarzmann (Eds.), Vol. 11704. 166–182.
- [3] Ethan Buchman. 2016. *Tendermint: Byzantine Fault Tolerance in the Age of Blockchains*. Ph.D. Dissertation. The University of Guelph.
- [4] Vitalik Buterin. 2013. A Next-Generation Smart Contract and Decentralized Application Platform. <http://archive.fo/Sb4qa>.
- [5] Josiah Carlson. 2013. *Redis in Action*. Manning Shelter Island.
- [6] Benoit Chesneau. 2021. Gunicorn. <https://gunicorn.org/>.
- [7] CodeNotary. 2021. immudb. <https://codenotary.io/technologies/immudb/>.
- [8] Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. 2010. Benchmarking Cloud Serving Systems with YCSB. In *Proc. of 1st ACM Symposium on Cloud Computing*. 143–154.
- [9] Katherine Cox-Buday. 2017. *Concurrency in Go: Tools and Techniques for Developers* (1st ed.). O'Reilly Media, Inc.
- [10] Rasmus Dahlberg, Tobias Pulls, and Roel Peeters. 2016. Efficient Sparse Merkle Trees: Caching Strategies and Secure (Non-)Membership Proofs. *Cryptology ePrint Archive*, Report 2016/683.
- [11] Hung Dang, Tien Tuan Anh Dinh, Dumitrel Loghin, Ee-Chien Chang, Qian Lin, and Beng Chin Ooi. 2019. Towards Scaling Blockchain Systems via Sharding. In *Proc. of 2019 ACM SIGMOD International Conference on Management of Data*. 123–140.
- [12] The PostgreSQL Global Development Group. 2021. PostgreSQL. <https://www.postgresql.org/>.
- [13] Dgraph. 2021. BadgerDB. <https://github.com/dgraph-io/badger>.
- [14] Tien Tuan Anh Dinh, Rui Liu, Meihui Zhang, Gang Chen, Beng Chin Ooi, and Ji Wang. 2018. Untangling Blockchain: A Data Processing View of Blockchain Systems. *IEEE Transactions on Knowledge and Data Engineering* 30, 7 (2018), 1366–1385.
- [15] Tien Tuan Anh Dinh, Ji Wang, Gang Chen, Rui Liu, Beng Chin Ooi, and Kian-Lee Tan. 2017. BLOCKBENCH: A Framework for Analyzing Private Blockchains. In *Proc. of 2017 ACM SIGMOD International Conference on Management of Data*. 1085–1100.
- [16] Muhammad El-Hindi, Carsten Binnig, Arvind Arasu, Donald Kossmann, and Ravi Ramamurthy. 2019. Blockchaindb: A Shared Database On Blockchains. *Proc. VLDB Endow.* 12, 11 (2019), 1597–1609.
- [17] Apache Software Foundation. 2017. Kafka. <https://kafka.apache.org/>.
- [18] The Linux Foundation. 2021. Hyperledger Fabric. <https://www.hyperledger.org/use/fabric>.
- [19] Johannes Gehrke, Lindsay Allen, Panagiotis Antonopoulos, Arvind Arasu, Joachim Hammer, James Hunter, Raghav Kaushik, Donald Kossmann, Ravi Ramamurthy, Srinath T. V. Setty, Jakub Szymbaszek, Alexander van Renen, Jonathan Lee, and Ramarathnam Venkatesan. 2019. Veritas: Shared Verifiable Databases and Tables in the Cloud. In *Proc. of 9th Biennial Conference on Innovative Data Systems Research*.
- [20] BigchainDB GmbH. 2018. *BigchainDB 2.0 The Blockchain Database*. Technical Report.
- [21] Miguel Grinberg. 2018. *Flask Web development: Developing Web Applications with Python*. O'Reilly Media.
- [22] Mike Hearn and Richard Gendal Brown. 2019. Corda: A Distributed Ledger. <https://bit.ly/3iLajrL>.
- [23] Dongxu Huang, Qi Liu, Qiu Cui, Zhuhe Fang, Xiaoyu Ma, Fei Xu, Li Shen, Liu Tang, Yuxing Zhou, Menglong Huang, et al. 2020. TiDB: a Raft-based HTAP database. *Proc. VLDB Endow.* 13, 12 (2020), 3072–3084.
- [24] MongoDB Inc. 2021. MongoDB. <https://www.mongodb.com/>.
- [25] Apache Kafka. 2021. Documentation. <https://archive.fo/tYpo6>.
- [26] Eleftherios Kokoris-Kogias, Enis Ceyhan Alp, Linus Gasser, Philipp Jovanovic, Ewa Syta, and Bryan Ford. 2020. CALYPSO: Private Data Management for Decentralized Ledgers. *Proc. VLDB Endow.* 14, 4 (2020), 586–599.
- [27] Jay Kreps. 2014. Benchmarking Apache Kafka: 2 Million Writes Per Second (On Three Cheap Machines). <https://archive.fo/8yUbt>.
- [28] Redis Labs. 2021. Redis. <https://redis.io/topics/benchmarks>.
- [29] Peilun Li, Guosai Wang, Xiaoqi Chen, Fan Long, and Wei Xu. 2020. Gosig: A Scalable and High-Performance Byzantine Consensus for Consortium Blockchains. In *Proc. of 11th ACM Symposium on Cloud Computing*. 223–237.
- [30] Ling Liu and M. Tamer Özsu (Eds.). 2009. *Merkle Trees*. Springer US, 1714–1715.
- [31] Senthil Nathan, Chander Govindarajan, Adarsh Saraf, Manish Sethi, and Praveen Jayachandran. 2019. Blockchain Meets Database: Design And Implementation Of A Blockchain Relational Database. *Proc. VLDB Endow.* 12, 11 (2019), 1539–1552.
- [32] Daniel Peng and Frank Dabek. 2010. Large-scale Incremental Processing Using Distributed Transactions and Notifications. In *Proc. of 9th USENIX Symposium on Operating Systems Design and Implementation*. 251–264.
- [33] Yanqing Peng, Min Du, Feifei Li, Raymond Cheng, and Dawn Song. 2020. Falcondb: Blockchain-Based Collaborative Database. In *Proc. of 2020 ACM SIGMOD International Conference on Management of Data*. 637–652.
- [34] PingCAP. 2021. TiDB. <https://github.com/pingcap/tidb>.
- [35] RedBeardLab. 2018. RedisSQL. <https://redisql.com/>.
- [36] Pingcheng Ruan, Tien Tuan Anh Dinh, Dumitrel Loghin, Meihui Zhang, Gang Chen, Qian Lin, and Beng Chin Ooi. 2021. Blockchains vs. Distributed Databases: Dichotomy and Fusion. In *Proc. of 2021 ACM SIGMOD International Conference on Management of Data*. 1–14.
- [37] Pingcheng Ruan, Dumitrel Loghin, Quang-Trung Ta, Meihui Zhang, Gang Chen, and Beng Chin Ooi. 2020. A Transactional Perspective on Execute-Order-Validate Blockchains. In *Proc. of 2020 ACM SIGMOD International Conference on Management of Data*. 543–557.
- [38] Felix Martin Schuhknecht, Ankur Sharma, Jens Dittrich, and Divya Agrawal. 2019. ChainifyDB: How to Blockchainify any Data Management System. *CoRR abs/1912.04820* (2019). arXiv:1912.04820
- [39] Ankur Sharma, Felix Martin Schuhknecht, Divya Agrawal, and Jens Dittrich. 2019. Blurring the Lines between Blockchains and Database Systems: The Case of Hyperledger Fabric. In *Proc. of 2019 ACM SIGMOD International Conference on Management of Data*. 105–122.
- [40] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. 2010. The Hadoop Distributed File System. In *Proc. of the 2010 IEEE 26th Symposium on Mass Storage Systems and Technologies*. 1–10.
- [41] Facebook Open Source. 2021. RocksDB. <https://rocksdb.org/>.
- [42] Tendermint. 2018. Documentation - ADR 020: Limiting txs size inside a block. <https://archive.fo/wip/o8Mhq>.
- [43] Jason Teutsch, Michael Straka, and Dan Boneh. 2019. Retrofitting A Two-way Peg Between Blockchains. *CoRR abs/1908.03999* (2019). arXiv:1908.03999
- [44] Sheng Wang, Tien Tuan Anh Dinh, Qian Lin, Zhongle Xie, Meihui Zhang, Qingchao Cai, Gang Chen, Beng Chin Ooi, and Pingcheng Ruan. 2018. Forkbase: An Efficient Storage Engine for Blockchain and Forkable Applications. *Proc. VLDB Endow.* 11, 10 (2018), 1137–1150.
- [45] Xinying Yang, Yuan Zhang, Sheng Wang, Benquan Yu, Feifei Li, Yize Li, and Wenyan Yan. 2020. LedgerDB: A Centralized Ledger Database for Universal Audit and Verification. *Proc. VLDB Endow.* 13, 12 (2020), 3138–3151.
- [46] Maofan Yin, Dahlia Malkhi, Michael K. Reiter, Guy Golan Gueta, and Ittai Abraham. 2019. HotStuff: BFT Consensus with Linearity and Responsiveness. In *Proc. of 2019 ACM Symposium on Principles of Distributed Computing*. 347–356.
- [47] Cong Yue, Zhongle Xie, Meihui Zhang, Gang Chen, Beng Chin Ooi, Sheng Wang, and Xiaokui Xiao. 2020. Analysis of Indexing Structures for Immutable Data. In *Proc. of 2020 ACM SIGMOD International Conference on Management of Data*. 925–935.
- [48] Meihui Zhang, Zhongle Xie, Cong Yue, and Ziyue Zhong. 2020. Spitz: A Verifiable Database System. *Proc. VLDB Endow.* 13, 12 (2020), 3449–3460.